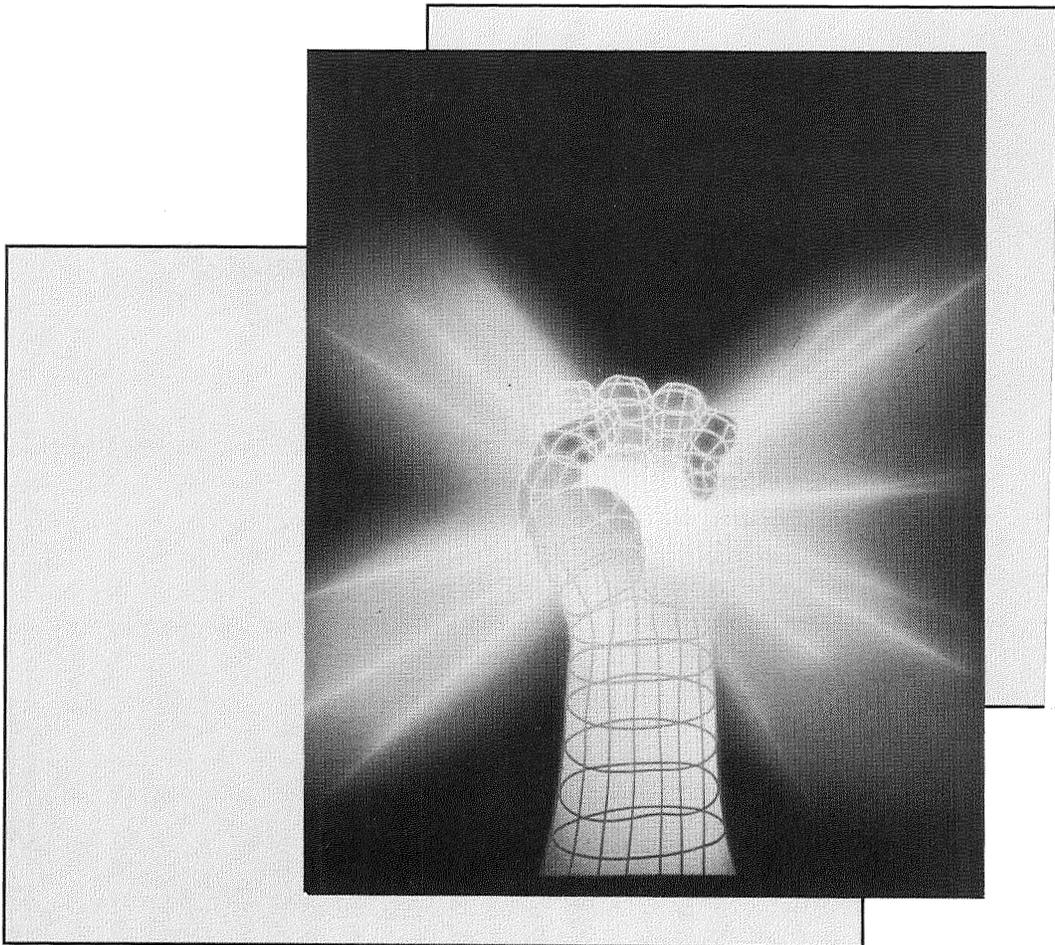


NASA Conference Publication 3189, Volume 1

The Third National Technology Transfer Conference & Exposition

December 1-3, 1992 • Baltimore, MD



N93-25561
--THRU--
N93-25617
Unclass

H1/99 0150470

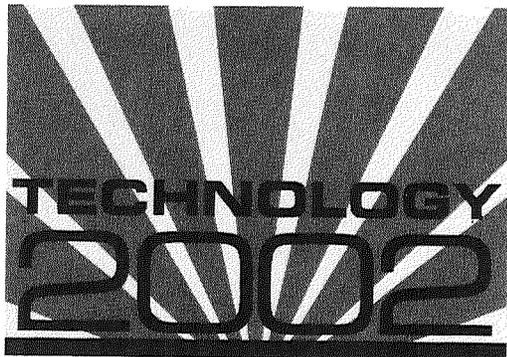
(NASA-CP-3189-Vol-1) TECHNOLOGY
2002: THE THIRD NATIONAL TECHNOLOGY
TRANSFER CONFERENCE AND EXPOSITION,
VOLUME 1 (NASA) 524 p

Conference Proceedings

Sponsored by NASA, the Technology Utilization Foundation,
and NASA Tech Briefs Magazine



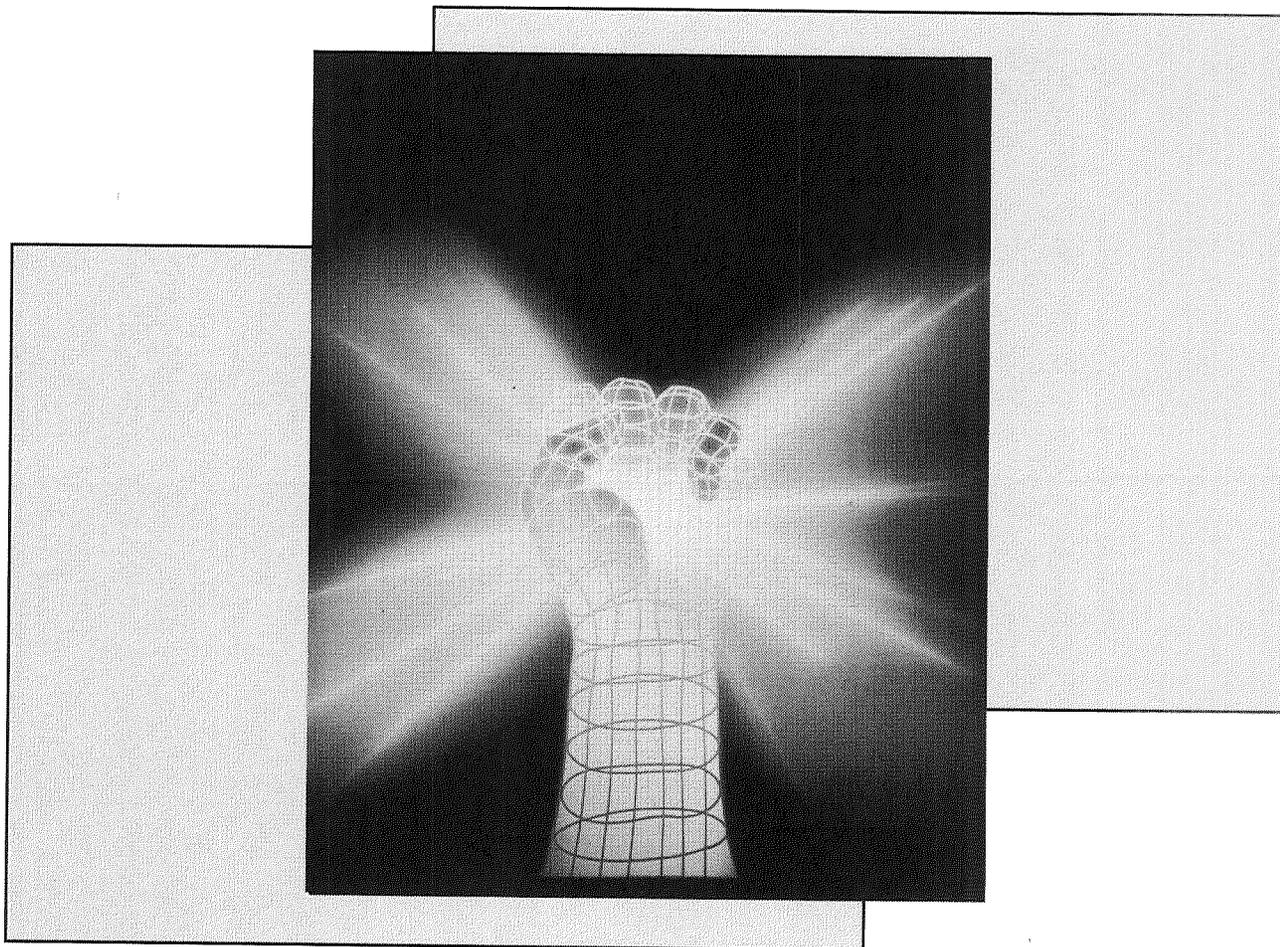
STI PROGRAM
SCIENTIFIC &
TECHNICAL
INFORMATION



NASA Conference Publication 3189, Volume 1

*The Third National Technology
Transfer Conference & Exposition*

December 1-3, 1992 • Baltimore, MD



Conference Proceedings

*Sponsored by NASA, the Technology Utilization Foundation,
and NASA Tech Briefs Magazine*



L.

FOREWORD

Technology 2002, the third national technology transfer conference and exposition, was held at the Baltimore, Maryland, Convention Center December 1-3, 1992. Opening ceremonies on November 30th, which featured remarks by Maryland Governor Donald Schaefer and a pre-conference reception, officially kicked off "National Technology Transfer Week." Technology 2002 shared the spotlight with the President's National Technology Initiative (December 1) and Massachusetts Institute of Technology's Entrepreneurial Technology Transfer Conference (December 3-5), highlighting the most comprehensive series of technology transfer events ever.

As Federal technology transfer becomes ever more important to our Nation's economic prosperity, the Technology 2000 series of conferences continues to lead the way, bringing together top leaders from government, industry and academia as it addresses issues vital to the enhancement of our competitive profile in the global marketplace. Welcoming remarks by NASA Administrator Daniel Goldin and Senator Barbara Mikulski's (D-MD) keynote address during Tuesday morning's opening agenda further accentuated this point.

This year's symposia reflects the importance of cutting-edge technology as a means to this end. Abstracts submitted from the R&D facilities of the Federal government and its contractors, addressing technological areas identified as critical by the White House, produced the most competitive field of technical presentations yet. But the development of the technology is only part of the solution. Forums on University technology transfer opportunities and foreign inventions available for U.S. benefit were conducted, along with briefings detailing existing Federal programs and opportunities currently available for use by industry, emphasizing the importance of establishing efficient mechanisms for the effective transfer of these technologies.

We are pleased to provide the Conference Proceedings from this year's symposia, during which 120 papers were scheduled for presentation. The Proceedings have been published in two volumes. Volume One contains the first sixty (60) papers scheduled for presentation, in order of presentation, while Volume Two contains the second sixty (60) papers, again in order. In addition, Volume One contains the complete transcript of Goldin's welcome and Senator Mikulski's address.

For information regarding additional sets of the Proceedings or audiocassettes of each presentation, please contact: The Technology Utilization Foundation, 41 East 42nd Street, Suite 921, New York, NY 10017, Ph.: 212/490-3999; Fax: 212/986-7864.

Sincerely,



Leonard A. Ault
Technology 2002
Conference Program Chairman

1-11
PRECEDING PAGE BLANK NOT FILMED

111
111

TABLE OF CONTENTS

National Technology Transfer Week	10
Opening Session	20
Keynote Address	50
 Symposia:	
Advanced Materials Part 1: Materials Processing	
INEL Spray Forming Research	15-1
Film Fabrication Technologies at NREL	25-2
The Effect of Hydrogen on the Optical and Scratch-Resistant Properties of Diamondlike Carbon Films	30-3
The Effect of Extrusion on PS-212 Self-Lubricating Materials	41-4
 Biotechnology and Life Sciences Part 1	
Measuring the Metastatic Potential of Cancer Cells	59-5
Immunoconjugates: Magic Bullets for Cancer Therapy?	71-6
Automated System for Early Breast Cancer Detection	77-7
Design of Mechanically Compatible Fasteners for Human Mandible Reconstruction	86-8
 Energy and Environment Part 1: Environmental Technologies	
Cone Penetrometer Measures Spectral Characteristics of Soils <i>In Situ</i>	99-9
Soil Reclamation and Recovery of Radionuclides and Toxic Metals	109-10
Converting Environmental Wastes into Valuable Resources	118-11
Low-Cost Dewatering Waste Slurries	121-12
 Information and Communications Part 1: High-Performance Computing and Networking	
High-Performance Networks and Supercomputers for Real-Time Flight Simulation	129-13
The SPLASH II Attached Processor System	139-14
Object-Oriented Tools for Distributed Computing	146-16
The Database Query Support Processor	156-16
 Manufacturing Technology Part 1	
"On Machine Tool" 3D Laser Measurement System	169-0
Application of an On-Machine Gauge for Diameter Measurement	170-17
On-Machine Capacitance Dimensional and Surface Profile Measurement System	178-18
Ultrasonic Polishing	182-19
 Microelectronics/Optoelectronics Part 1	
Two- and Three-Dimensional High-Performance, Patterned Overly Multi-Chip Module Technology	195-20
Improved Performance and Safety for High-Energy Batteries	205-21
Thin Rechargeable Batteries for CMOS SRAM Memory Protection	213-22
Passive Stacking for Improved Vibration Isolation	219-0
 Advanced Materials Part 2: Ceramics and Composites	
A Novel Method for Characterization of Superconductors: Physical Measurements and Modeling of Thin Films	223-23
Production of Ultrafine, High-Purity Ceramic Powders	231-24
Mullite Whiskers and Mullite-Whisker Felt	241-26
Graphite/Epoxy Composite Laminates with Co-Cured Interlaminar Damping Layers	250-26

Artificial Intelligence Part 1

Expert System for UNIX System Reliability and Availability Enhancement	259	-27
The Generic Spacecraft Analyst Assistant: A Tool for Developing Graphical Expert Systems	268	-28
TARGET: Rapid Capture of Process Knowledge	279	-29
Tree Classification Software	289	-30

Biotechnology and Life Sciences Part 2: Medical Technology

Automatic Detection of Seizures	301	-31
A Fiber Optic Probe for the Detection of Cataracts	308	-32
Heart Rate Spectral Analysis System	317	-33
The Application of Integrated Knowledge-Based Systems for the Biomedical Risk Assessment Intelligent Network (BRAIN)	322	-34

Energy and Environment Part 2: Energy Innovations

Solid-State Isotopic Power Source for Computer Memory Chips	335	-35
Photovoltaic Power Without Batteries for Continuous Cathodic Protection	341	-36
High-Speed Solid-State Circuit Breaker	345	-37
Variable-Speed Generators with Flux Weakening	353	-38

Information and Communications Part 2: Computer Simulation and Modeling

Industrial Applications of Computational Fluid Dynamics	365	-0
Scientific Visualization Using the Flow Analysis Software Toolkit	366	-39
Integration of Design, Thermal, Structural, and Optical Analysis, Including Thermal Animation	376	-40
Data Systems Dynamic Simulator	385	-41

Manufacturing Technology Part 2

A Novel Optical/Digital Processing System for Pattern Recognition	397	-42
Vision-Aided Monitoring and Control of Thermal Spray, Spray Forming, and Welding Processes	407	-43
Automatic Robotic Variable Polarity Plasma Arc (VPPA) Welding	415	-44
Firmware Development Improves System Efficiency	425	-45

Advanced Materials Part 3: Plastics, Polymers, and Rubbers

Electro-Expulsive Separation System	437	-0
Elastomer Compound for High-Wear Applications	438	-46
Dynamic Hardness Tester and Cure Meter	446	-47
Instrumentation Measures Gas Permeability of Polymeric Membranes	451	-48

Artificial Intelligence Part 2

A Software Package for Neural Network Applications	463	-49
Control of Complex Dynamic Systems by Neural Networks	473	-50
Adaptive Process Control with Fuzzy Logic and Genetic Algorithms	483	-51
A Genetic Algorithm Tool for Complex Scheduling Problems	491	-52

Energy and Environment Part 3: Environmental Technologies

Development of a LIDAR Mapping Instrument	503	-53
Commercial Applications of a Multispectral Sensor System	515	-54
Interactive Forecasting with the National Weather Service River Forecast System	527	-55
Application of Space Life Support Technology to Terrestrial Environmental Problems	537	-56

PRIMARY

NATIONAL TECHNOLOGY TRANSFER WEEK

**National Technology Transfer Week
November 30 - December 5, 1992**

"Working To Keep America Strong"

Three high-powered events that came together to help U.S. Industry utilize cutting-edge technology to gain a competitive edge in the international marketplace.

**Technology 2002, The Third National Technology Transfer Conference and Exposition
December 1-3, 1992 at the Baltimore, Maryland, Convention Center**

America's premier technology transfer showcase, Technology 2002 featured 120 presentations spotlighting inventions with commercial potential in manufacturing, materials, computing, communications, biotechnology, and energy/environment, areas identified by the White House as National Critical Technologies. Hands-on workshops on patent licensing, cooperative R&D, and research grants, as well as sessions on university-based and international technologies were also conducted. Over 60,000 square feet of exhibits showcased the latest products and processes available for license or sale.

**The President's National Technology Initiative (NTI)
December 1, 1992 at the Baltimore, Maryland, Convention Center**

Held concurrently with Technology 2002's Tuesday sessions, the NTI focussed on opportunities for partnerships between government, academia, and U.S. companies to translate new technologies into marketable goods and services. Top-level officials from the White House, the departments of Commerce, Energy and Transportation, NASA, and other Federal agencies discussed three critical elements for U.S. competitiveness: technology, capital and manufacturing.

**Massachusetts Institute of Technology's Entrepreneurial Technology Transfer Conference:
The Commercialization Success Factors for Intra- and Entrepreneurs
December 2-5, 1992 at the Baltimore, Maryland, Hyatt Regency Hotel**

This world-class conference sponsored by the Massachusetts Institute of Technology Enterprise Forum, is designed to provide the skills and tools needed to commercialize emerging technologies and capitalize on partnership opportunities such as those featured at Technology 2002 and the NTI. In more than 25 "how to" sessions and interactive tutorials, leading entrepreneurs, technology licensing officers, research managers, and investors helped to formulate winning technology transfer strategies and techniques.

1.

OPENING SESSION

Conference Opening and Federal Technology Overview

After Monday evening's Opening Ceremonies, Technology 2002 activities got under way Tuesday morning, December 1st, with the Conference Opening and Federal Technology Overview. NASA Administrator Daniel Goldin officially welcomed conference attendees and Senator Barbara A. Mikulski (D-MD) delivered the morning's keynote address. The following is the complete transcript their remarks.

Introduction by Bill Schnirring:

Good Morning, Ladies and Gentlemen. Thank you for coming to the opening Plenary Session that kicks off National Technology Transfer Week, comprised of Technology 2002, the President's National Technology Initiative and the MIT Enterprise Forum. I'm Bill Schnirring, president of Associated Business Publications, the publishers of NASA Tech Briefs magazine, which is, along with NASA and the Technology Utilization Foundation, a co-sponsor of Technology 2002. The objective of Technology Week and Technology 2002 is to accelerate the transfer of technology from the public sector to the private sector, and to enhance America's competitive and ultimately benefit all mankind.

It is my privilege this morning to introduce you to a man who began his working career with NASA in 1962 as a research scientist. After five years at the Lewis Research Center, he moved into the private sector and began a long and distinguished career at TRW, where, until the May of this year, he was Vice President and General Manager of TRW's Space and Technology Group. On April 1st, 1992, he was named the ninth administrator of the National Aeronautics and Space Administration. Ladies and Gentlemen, please join me in a warm welcome for NASA Administrator, Daniel S. Goldin.

Daniel S. Goldin:

Thank you Bill. I'd like to compliment you on the ultimate in tech transfer. NASA, in 1985, turned over the responsibility to tech transfer to Bill and his organization at Tech Transfer Briefs. And at the time, the circulation was 70,000, and under Bill's guidance, without any government subsidy, the circulation is now up to 200,001. It's the second largest engineering publication in the world, and Bill, you're to be complimented on this wonderful task you've done.

One hundred and fifty years ago, the famous the French social observer, Alexis de Tocqueville, came to our shores. Came to this nation that literally invented itself to see what made it's people different with those from the old world. He wrote "America's a land of wonders in which everything is in constant motion, and every change seems an improvement. No natural boundaries seem to be set to the efforts of man, and in his eyes, what is not yet done, is only what he has not yet attempted to do."

De Tocqueville wrote this back in the age of keel boats and canals, an age that we would say moved pretty slow, we had to go west, out to the frontier. We had to build railroads, super highways, jet liners or space craft. But even back then, he knew what American's were capable of, that we could do whatever we attempted to do.

Today, America is still in constant motion. The question is for our nation, are we still bold? Are we still willing to take risks and attempt the seemingly impossible, or will fear cause us to seek a false sense of security and comfort by sticking with the safe and the familiar.

America has to get bold again. If we're going to shape our own destiny, rather than have the future shape us, we've got to take some risks and make significant investments. We've got to stop eating our technical seed corn and plant it, work the soil and then bring in the harvest. This is the two-hundred year tradition of America, the careful legacy inherited from every pioneer and inventor that came before us.

For everyone who's worried about the American economy, being stuck in a rut, it's vital that we remember the tremendous power of technology to cause growth. During the days of Apollo, American reigned supreme in the world of technology because we were one the cutting edge of it in space. Over the decades, NASA has generated at least three, 30,000 known spin-offs from it's technology, creating new products in the industries worth billions that have changed the face of America.

But there has not really been enough of a focus, systematic commitment on NASA's part to the technology transfer process. For years, tech transfer has been the lonely foster child and that description could apply to other government agencies as well. America created a large federal research effort to fight the Cold War and we fought it very well, but now it's over. America's needs have changed, and the Federal Government must respond to those needs by putting America's Cold War technology base to work for our economy.

There is a peace dividend beyond just the reduced threat of war and lower defense budgets, as important as they are. This is the technology dividend, as we multiply the government's technology transfer efforts. That's why we have started the National Technology Transfer Initiative, thanks to the vision of Secretary James Watkins, to let the sun shine into the federal laboratories. That's why 50 federal labs from 13 agencies are here today, and that's why you're here, to find out what we have to offer.

NASA itself is changing with the times. To dramatically improve the way NASA approaches both the development and transfer of technology, and the commercialization of space, we have created the office of the Advanced Concepts and Technology. This will be an entirely new breed. A highly flexible, customer-driven organization that will develop innovative concepts and high leverage technology that both fulfills NASA's needs, and have very, very significant commercial possibilities.

One of the primary functions of this office is to be NASA's front door to the businesses who want NASA's help and expertise, who have new ideas and technologies for us. Our new office will provide one-stop shopping for technology, customers and suppliers, whether they are businesses or universities, or even program offices within NASA.

For example, say a small woman-owned firm that makes thermo plastics asks for our help in developing a new manufacturing technique. We could develop a partnership whereby the firm receives technical expertise that improves her product line, while NASA receives new lightweight materials for our own use at lower cost and space. Such a win/win approach would enhance NASA's programs, yield more value to the tax payer and improve the economy generally by helping the private sector become more competitive.

The Office of Advanced Concepts and Technology will also set up new mechanisms and improve existing mechanisms to aggressively transfer technology into America's economy. We will do this by seeking the input of technology user community to figure out the best transfer mechanisms, whether it's by reading, publications like Tech Briefs, regional tech transfer centers, centers for the commercial development of space, cooperative research agreements, or actually having you come and work in our laboratories for jointly developed products.

In fact, Greg Reck, the acting head of this new office, is holding meetings around the country to solicit inputs on how NASA can improve it's tech transfer process. I really urge you to contact this office if you have any ideas on how we could do the job better. But let me stress, customer need must govern the type of tech transfer mechanisms that we devise. If it's user friendly, it's not going to maximize the benefit for the American economy.

The true test of NASA as a jobs generator is not how many people are working for NASA, but how many people are working in America because of NASA, because of NASA's ability to reach out into the future and bring the answers back for today.

Just this year, a NASA-derived laser was approved by doctors to use a cleaning, approved for doctors to use in cleaning out plaque from blood vessels. In many cases, this simple \$8,000 procedure can replace a \$25,000 coronary by-pass operation, with much less risk and trauma to the patient. More than 100 people are now working in a brand new company that builds this device.

Wireless infrared head sets, invented for the space shuttle communications, have been commercialized and are revolutionizing personal communications on the factory floor, in concert halls for the hearing impaired, in language translation equipment, and can eliminate the need for computer network cabling inside the office.

High temperature material developed by NASA enabled an Ohio steel making equipment firm to triple a lifetime of very expensive rollers used in the continuous casting of steel, resulting in substantial sales increases and improving the efficiency of American steel plants all over the country. Even the paint protecting the Statue of Liberty, who holds a torch high over the land of the free, was developed by NASA.

I believe NASA can be a leading force in our society. The discoveries and technology the space program provide inspiration for our minds and souls, and hope for our future survival, opportunity for renewed prosperity, and catalytic action for peace through international partnerships. NASA can do all this and more, if America continues to be bold and keeps reaching for the stars.

In the early part of the next century, a passenger jet will take off from California and land in Tokyo in a little over four hours. A child will sit at a classroom computer and use a virtual reality terminal to explore and understand the interior of the human body. And a few decades later, two explorers will set out from their base camp and look for subsurface water, mineral resources and evidence of life on Mars. I want that plane to be built in America, I want that child to be an American child, and I want those astronauts to wear the American flag right here on their shoulder.

The challenge of developing technology and reaching for things we can't yet grasp is what made this country great, and will keep us forward as long as we have the courage to live on the cutting edge. Thank you very much.

KEYNOTE ADDRESS

Introduction by Daniel S. Goldin:

I now have the great pleasure of introducing our keynote speaker, a Baltimore native who for years has been hard-working and steadfast servant to the people of Maryland. Senator Barbara Mikulski has championed the space program because she has a deep and profound understanding of the critical importance of research and technology to the future of this country and its economy. As chairman of the appropriations subcommittee that handles NASA budget, NASA's budget, Senator Mikulski has been a staunch supporter of the Space Station Freedom, because she understands not just its critical role as a place where we could learn how to live and work in space, but as the first major U.S. facility in permanent orbit, it will enable cutting-edge research into materials, biomedical and life science that will pay enormous dividends here on Earth, generating new technologies, new industries and new jobs.

Through Senator Mikulski's visionary leadership, this year, NASA and the National Institutes of Health, signed a new cooperative agreement, so that we can both maximize our unique resources in life science research, both on the ground and on space station freedom. We can also work together on medical technology development.

And just one effort we're excited about, our technology transfer specialists are working with the National Cancer Institute to apply advance digital imaging to mammography, which could really improve the early detection of breast cancer in women.

NASA and NIH can achieve great things together to improve the health of all Americans, and Senator Mikulski, we deeply appreciate your strong support in making it happen.

Senator Mikulski was also visionary in strongly urging the United States and NASA to begin working with Russia in a way of building bonds that lead to a lasting peace. I followed up on her urgings, and NASA went to Russia, and we concluded a negotiation to work with the Russians in space. And as a result, in November of 1993, a cosmonaut will fly in the shuttle, and after that, in 1995, one of astronauts will fly on board the MIR Space Station, where we will conduct life science experiments, and then the shuttle will rendezvous with MIR with an advanced medical laboratory inside, to do fundamental life science research in space.

And it's an exciting way to start a relationship that we hope will result in a stable peace and scientific benefits for all human kind. And once, again, Senator Mikulski. It's been a super experience working with you this year as the NASA Administrator, and it is my great honor to introduce to you all Senator Barbara Mikulski.

Senator Barbara A. Mikulski:

Good Morning everybody. Boy, you don't like talking in a dark room like this? First of all, I can't see you, I guess you can see me. Why do we have the lights down? Is there any reason for that? I hope that you're so mesmerized by the charismatic speech that Dan Goldin just gave, the magnetism that I'm about to project, the wisdom of Admiral Watkins, the can-do attitude of Barbara Franklin, that you're going to be taking notes in an obsessive and compulsive way. I see that the lights are up, and you know, this is kind of a new morning in America here. Thank you.

I know the lights went down so we could have that stirring rendition of the National Anthem, which of course, was written here in my home town of Baltimore.

I want to welcome all of you to this great city, and this town that I call my own hometown. We feel that we in Baltimore have several major league teams. One is called the Baltimore Orioles, and we invite you back this summer. But we have other major league teams here in the town Baltimore and in the state of Maryland. Those teams are called the University of Maryland. They're called the Johns Hopkins University. They're called federal laboratories, like Goddard, NIH, NIST, FDA, all of which are right here in our own state of Maryland.

And what we feel are really the major, are the infrastructure, the intellectual infrastructure, to make sure the United States of America is in the major leagues from now and well in to the twenty first century.

I'm so pleased that you've chosen this City to hold this Technology Transfer meeting. And I also am very pleased and honored to be on the platform with the distinguished people that you see here today, representing Cabinet level participation.

I want to express a particular thanks to the Bush Administration for the way that they are working and cooperating as we transit from one party to another, in terms of the Executive Branch. The mandates that have been established by President Bush that this transition will be not only orderly, not only timely, but will trans, the transfer will occur in a way for which we not miss a beat in what we're doing, is very much appreciated.

Each and every person that I've worked here, and I'll speak to the two Cabinet level people, has really indeed been an honor. Admiral Watkins and I first met when he took on the difficult chore of chairing the AIDS Commission. And then we have worked together not only at the Department of Energy, but really his great passion in life is to make sure the Federal Laboratories are opened up in a way for both technology transfer, but also he sees the Federal laboratories working with local school systems.

Barbara Franklin I think has been an outstanding Secretary of Commerce in her short time, and she has brought back in the Department of Commerce just a great spirit of "Can-doism". And Dan Goldin and I have made, I think, a wonderful working relationship and I look forward to working with him on a continued basis.

A lot is happening how, as we look at where we are and where we've been. And I've very pleased that J.R. invited me to come and speak to the, to this transfer tech. I met J.R. a little over four years ago, when I became the subcommittee chair of V.A., HUD and Independent Agencies, in the Senate Appropriations Committee. The Senate Appropriations Committee, as you know, is the subcommittee, or the committee, within the United States Senate, that actually puts money in the Federal checkbook.

Now you have to know I love being in a United States Senator, that's why I work my earrings off to get re-elected, and ended up with seventy-one percent of the vote. But one of the problems that I have with being a United States Senator is that when all was said and done, more gets said, than gets done. And one of the things that I've loved about being on the appropriations committee, is that we are the only committee within the United States Senate and United States Congress that has a do-gate.

Every October 1st, we have to produce an appropriation, which means that we have to march in a very timely pace to accomplish that. When I became the Chair of V.A., HUD, and Independent Agencies and you might say how in the hell does that happen? Not how did I become chair, that was my own talent [laughter]. Modesty, of course, being a hallmark of it. That years ago, within the appropriations committee, independent agencies were created and they didn't know where to put them. So they created a special committee with Independent Agencies. They had names like NASA, HUD, V.A., but a lot of these agencies have now grown up to be Cabinet level responsibility. That's why I have the wide range of portfolio than just about anybody outside the Defense Subcommittee that Senator Inouye chairs.

And when I became the chair of that Subcommittee, I knew a lot about health care, and I knew a lot about Veterans. I knew a lot about housing, because I work in it as a social worker and as a City Council woman in this great city, and it has been something that I've worked with for a substantial part of my public career.

But the whole area of science and technology, and what it means to America's future, was in many ways new to me. And it was being able to get out, travel with NASA, its leadership and moving out to the space centers, as well as meeting with Eric Block at the National Science Foundation, that I really cut my teeth in terms of the possibilities of what this means to generate jobs today, and jobs tomorrow.

One of the people who really guided me in my early days on that Subcommittee was J.R. Thompson down at the Marshall Space Center, which was a tremendous help to me to understand the Space Station, it's role and it's mission, and both space, and what the Space Program has meant to the manufacturing in the United States of America.

So to you, J.R., I'm going to thank you for all that you've done for me, and what you've done for the United States of America. Wherever you are, I don't see you up here at the--there you are. I think we ought to give J.R. a wonderful round.

When I graduated from the Institute of Notre Dame back in the '50s, and went to Mount St. Agnes College, I thought maybe I would, I started out in pre-med. I thought I might be another Madame Curie. I thought I'd never be a United States Senator. In 1954 the politics in Maryland and Senators had pot-bellies, smoked cigars and certainly didn't look like or sound like Barb Mikulski.

So my election to the United States Senate and the election of four new democratic women to the United States Senate I think is a symbol of the new world order. And this new world order can be either a great tragedy for the last half of the 21st Century, or the end of the Cold War could be an end, or beginning of a cornucopia of opportunity.

There are those who look at the Cold War and see enormous possibility. I happen to share that vision. And happen to believe that where our greatest strategic threats are not military but are economic. And where our economic strength will largely rest in our successful development of technology, the transfer of technology and making sure that our young people have the education and skills to put that technology to productive use. And where we face the tremendous task of helping the corporations and the men and women who helped win the Cold War over the last five decades, make the transition to a civilian economy. The country that says that it wants to compete is the country that will make sure that it will lead.

The election of Bill Clinton and Al Gore on November 3rd indicated that our country was ready to face the challenges of the new world order. I believe that the Clinton/Gore administration, working with the new Congress, will usher in a new era for real partnerships between the Federal Government and the private sector in science and in technology.

I believe that with this election that the pledge that we who represent the Democratic Party need to make is that gridlock is over, stagnation is dead, and that we need to work in the United States Congress with an Executive Branch, that really makes sure that government is not part of the problem, part of the obstacles in developing ideas, technology and technology transfer, but actually is a source of being able to bring about the invention of technology, the transfer of technology and always changing the culture of the organizations involved with it.

I do not believe that the Clinton/Gore Administration will want to be involved in picking winners and losers among companies or ideas, nor in the development of five-point plans like you'd expect in some Trotsky pamphlet in the collapsing socialist empires around the world.

I believe that they want to make the government more entrepreneurial, more responsive to day-to-day needs and to use the spirit of innovation and flexibility that has been our country's hallmark, to reinvent the role of government. I do not believe that Bill Clinton and Al Gore want to make government bigger, nor spend more money, but make sure it works better and spends those precious tax dollars to make sure that we have a greater outcome.

I believe that it is their approach to make federal investments in research and development, but do that in our federal laboratories that will generate the pre-competitive ideas and technologies and have the private sector value-add in the product development.

Let me give you a few examples of where we know that technology transfer has worked just like that. Dan Goldin gave you some of the examples this morning. We all know the success story of the insulin implant for young diabetics. Or robots with the capacity to handle tasks too dangerous or difficult to do that. Those are the kinds of things we've done, and we'll be doing more.

But what I'd like to spend my time with you this morning is really to share with you what I think the next year, or the next four years will look like. I am not here to be Bill Clinton by proxy. But I've worked with those guys, and I particularly have worked with Senator Al Gore. Many people don't know that Senator Gore came to the House of Representatives together. This year, fifteen years ago, in 1976. We were part of the bicentennial class, and we both won coveted seats on the Energy and Commerce Committee. For eight years, I sat next to Al Gore as we worked on Energy policy, as we worked on telecommunications policy and a whole variety of other things.

During those long hearings, and difficult mark ups, when a committee chaired by John Dingle, you get to know your seat mate pretty well. And for eight years, Al Gore and I sat together trading notes, trading ideas, and trading strategies.

When we both came to the United States Senate, he a few years before me, we then had a new relationship, because Senator Gore chaired the authorizing committee on Space and I chaired the appropriations committee. And in that time we developed a lot of ideas, essentially on how we think government should work.

And, in a nutshell, it is our belief that we need a new technology policy for a new world order. A technology policy that generates smart kids, smart workers, smart managers and a smart government that works with a smart private sector that helps bring about this new world order.

One of our concerns over the years has been that federal agencies did not cooperate with each other, and for many years, we have talked about changing the culture of these agencies. It is, was a source of great concern to us, and I know it's shared by Dan Goldin, that very often it is the very culture of our institutions that in fact impede our development.

Many government agencies are based, as is our private sector, in a large, hierarchical, top-down, trickle down idea approach, whether it's General Motors, whether it's the Pentagon, or even, Dan, you would agree, NASA itself.

And what we find, just as the private sector is changing, so does government have to change. And just like it is incumbent for Congress to make sure working with the President that gridlock is over, we would encourage those people working in Federal agencies to stop worrying about turf, fiefdoms, and power struts and battles, and start focussing on the mission. Too must, just as within the Congress itself, we have too many committees, too many subcommittees, too much process, little outcome, and with staffs who continually jockey in their wingtips for greater and greater power accumulation, we see that also in the Federal agency.

I want to be part of the Congress of the United States that trims (sic) down and streamlines the working of the Congress, and when we talk about working with our Federal agencies and our Federal laboratories, we need to start cooperating with each other and let's stop worrying about who is going to be the leader in science education, who is going to be the leader and get the title for leading the way in high tech computer communication. Let's start worrying about how the United States of America is going start being the leader.

Let's learn from the Japanese. The society that makes a culture of cooperation is the society that competes the best. We have a culture of competition and we end up losing ground, losing money and losing time. So it is the spirit of cooperation that I think needs to be a hallmark of what the rest of this century's going to look like.

I happen to also believe that the President of the United States who has now said he wants to elevate the economic security at the same level as the national security, does well to look at the same framework for science and technology. You cannot have economic growth in the United States of America unless you have a science, technology manufacturing base and therein lies the key.

I hope that the President of the United States will turn to Al Gore as Vice President, and make sure that he is the top gun to help develop the administration's technology policy and help be that arbitrator that brokers the conflict, brokers the turf war, and then provides the leadership that will mean.

But let me then just say a few points about what I think this whole new administration will look like. I talked about the end of gridlock, I talked about the end of stagnation, but let me say something about presidential leadership. A lot of people are saying, "well, who is Bill Clinton? What's he going to do? What's he like, what does he think?"

Well, read the book "Putting People First". I think it will give you a lot of clues. But also look now, because as technology transfer people, you believe in outcomes and not process. Look at what he has done already and therein lies the clue. And look at who he is turning to for advice. First of all, one of his first executive decisions as the Democratic nominee was to pick his Vice President. And he picked Al Gore. This is a president, who, or a democratic nominee, who was not shy about picking someone of his own generation, someone of talents equal, in a different type of experience, and had a team to help him select Al Gore.

When you take a look at what they're doing on the cluster teams, I would encourage you to follow those because they will give you insights into the direction of where we're going. I'm not talking about reading tea leaves, Washington's abuzz reading tea leaves. Everybody reads the Washington Post, and the New York Times -- "What's it mean? Who was at Pamela Harriman's? Were you on the A list (and by the way, I was)?" I made the Dean's list.

But when we take a look at it, let's take a look at the Cluster Teams, and if we look at also who's heading the economic strategy. If you read the works of Robert Reich, you will get a clear sense of where Bill Clinton really wants to lead the United States of America. Robert Reich talks about the development of technology, of technology transfer and making significant investments in young people, to make sure that we have a skilled educated work force, to not only develop the new technology but to operate that technology. And that government spending be organized not to simulate consumption but to stimulate growth and to make long term sustainable commitments in what we're going to do.

And I know for those of you in government, one of the concerns that you have is that with every damn appropriations, you never know if your project's going to be here today or gone tomorrow. And I share that frustration with you. Also the fact that Bill Clinton and the Transition Team has asked Dr. Sally Ride to oversee what they're doing in science and technology gives you an insight in that.

I would hope that what the Clinton Administration does then is to move to real teamwork among Federal Agencies so that we could have real outcomes. We want to make sure that as they then move ahead, that the other principle that they follow is to undertake a review of Federal Tax and Regulatory policy to create a climate and investment in science technology. We want progrowth, that means reinstating the investment tax credit and making permanent the research and development tax credit for our private sector. And at the same time, I believe, we need to take a look at our anti-trust laws that were created for 19th century monopolies, rather than for a 21st century economy, that actually impedes the alliances and consortiums that can be formed between government, the private sector, and our universities. And I also want to make sure that if you invent an idea, you get to keep it, and I believe that one of the strongest things that we need to do, both in our country, and in our trade negotiations, is to make sure we protect American intellectual property.

Those are the kinds of things that we think we need to be doing. As well as to bring in U.S. industry and America's research universities together to strengthen our competitiveness. You know that a recent Business Week story highlighted those areas to be considered the new "Silicone Valleys" and what we need to do to be competitive.

And what do they say makes it? What we need is visionary entrepreneurs, world-class universities and the presence of a government at all levels that is a catalyst and a facilitator, not as a super board of directors of industry.

These are some of the things that we hope change. But it does mean also breaking the cultural mind-set that shuns the transfer of technology and knowledge in selecting what we do. Too often the culture of science and technology in the Federal laboratories is to ask for more and more money for individual investigators without a clear, national navigational chart on where we're going. I support although the premise of all those wonderful investigators out there who have more ideas than ever could be funded.

But the United States of America no longer can fund good intentions, no longer can fund every good idea to be pursued. We need a national navigational chart on where we're going for the 21st century, and I believe if we focus on the development of critical technologies, this will enable all those young investigators to have all the opportunities to pursue research, but at the same time the culture needs to change not only from the pursuit of ideas, but to the development of those ideas into technology and a technology into product development. Our United States of America continues to win Nobel Prizes while we continue to lose the markets. We hope that changes.

These are some of the things that I hope the, and I think the Clinton/Gore Administration will do. There are many other things that we'll be talking about in the days ahead. But the hallmarks of what we see as change is really the end of gridlock, the end of stagnation, the change of our culture so that we emphasize competition and that all federal spending will have to meet the test, when how does this generate a job today or a job tomorrow.

Dan Goldin said it best in his opening remarks when he said we should be measured not by how many people work for NASA but how many people in the United States of America work because of NASA? And how about because of the National Institutes of Health? How about because of the Federal Drug Administration? How about focusing our orientation and our thinking in that line.

I hope that when I join you next year at the Technology Transfer Conference, we can look back over a year where Congress has got it's act together, worked with an Executive Branch, worked with all of you who have been trying so hard, working so hard in the past, but wondering where were we going as a nation? And then to look back and say, we've made not only a good start, but we've jump started this economy. Like all of you, I want to make sure that when this century ends, we will be on the brink of where the best days of the United States of America are ahead of it, and not behind it.

I believe the Clinton/Gore team will do that. I think the legacy of George Bush needs to always be acknowledged. Because it was George Bush and the GI generation that enabled the next generation of baby boomers to come along. If it had not been for the GI generation, and gallant men like George Bush, who went to war to save Western civilization from it's threats, we would never have developed the opportunities and the society that we have now.

But now there is a turning of the wheel. And we want to salute those who have gone ahead of us, and created a framework for the United States of America, who have worked to bring about peace in the world, who by their very efforts and leadership have brought the end of the Cold War, that now enable us to have the opportunity to fight the new war, the new war for America's future. And I hope we can do it with the honor, and the steadfastness and the same ability that those have gone before. And I think we can do it and we rely upon you, the technology warriors of the future.

Thank you very much and I look forward to working with you.

omit

**ADVANCED MATERIALS PART 1:
MATERIALS PROCESSING**

PRECEDING PAGE BLANK NOT FILMED

~~PAGE 12~~ INTENTIONALLY BLANK

13.

INEL SPRAY-FORMING RESEARCH

5, -27
150471

Kevin M. McHugh
Idaho National Engineering Laboratory
Idaho Falls, ID 83415-2050

NO 3-25562

James F. Key
Idaho National Engineering Laboratory
Idaho Falls, ID 83415-2050

ABSTRACT

Spray forming is a near-net-shape fabrication technology in which a spray of finely atomized liquid droplets is deposited onto a suitably shaped substrate or mold to produce a coherent solid. The technology offers unique opportunities for simplifying materials processing without sacrificing, and oftentimes substantially improving, product quality. Spray forming can be performed with a wide range of metals and nonmetals, and offers property improvements resulting from rapid solidification (e.g. refined microstructures, extended solid solubilities and reduced segregation). Economic benefits result from process simplification and the elimination of unit operations. Researchers at the Idaho National Engineering Laboratory (INEL) are developing spray-forming technology for producing near-net-shape solids and coatings of a variety of metals, polymers, and composite materials. Results from several spray-forming programs are presented to illustrate the range of capabilities of the technique as well as the accompanying technical and economic benefits. Low-carbon steel strip >0.75 mm thick and polymer membranes for gas/gas and liquid/liquid separations that were spray formed are discussed; recent advances in spray forming molds, dies, and other tooling using low-melting-point metals are described.

LOW-CARBON STEEL STRIP

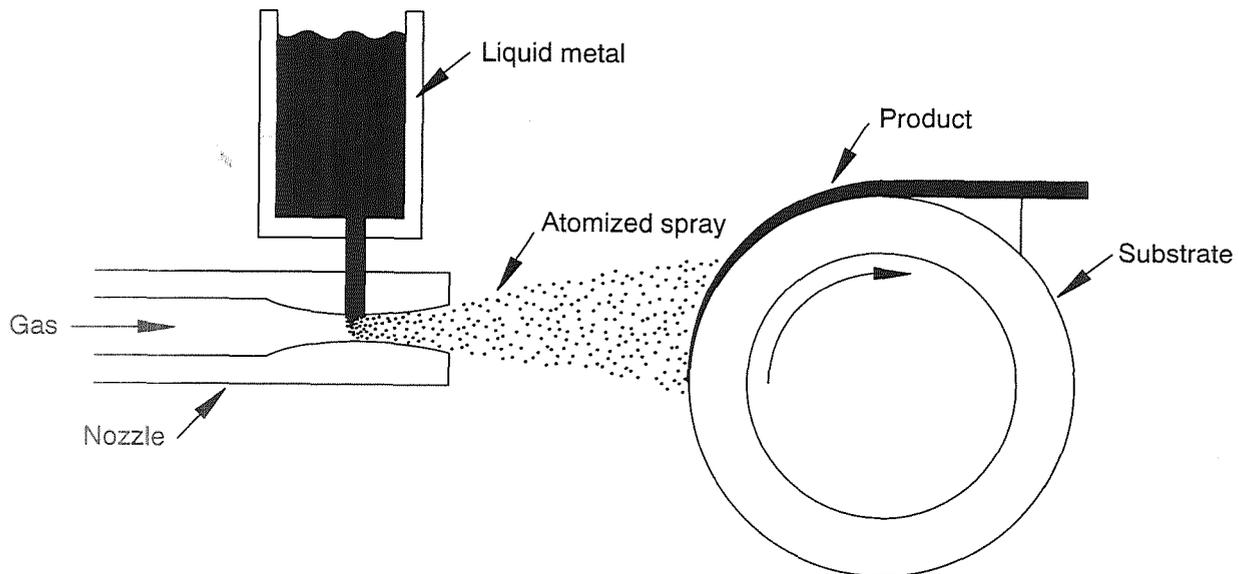
The National Critical Technologies Panel recently determined that materials processing is a leading critical technology for meeting future national needs [1]. Spray-forming technology may help address these needs by improving product quality while simultaneously lowering production costs. Nearly all low-carbon steel strip is produced by conventional ingot or thin-slab metallurgical techniques. The molten steel is cast as an ingot or slab and extensively deformed to obtain the desired shape and properties. This is highly energy intensive and requires large capital investments. In contrast, INEL spray-forming technology transforms molten metal to close to final strip form in a single rapid solidification step. Minor hot rolling then fully densifies and further refines the metal's microstructure. This can lead to enormous cost savings. For low-carbon steel hot band, the industry's highest volume commodity, technoeconomic analysis estimated improvement in production costs by as much as \$50 to 100 per ton. At the current domestic production rate of 50 to 60 million tons per year, this corresponds to cost savings of \$2 1/2 to 6 billion per year. This savings, due primarily to elimination of unit operations and associated energy costs, would give the United States a tremendous competitive advantage in low-carbon steel production. Improvements in product quality are equally impressive. For example, after hot rolling (~60% thickness reduction), INEL spray-formed low-carbon steel strip typically had about 50% higher yield and ultimate tensile strength than commercial strip, i.e., its properties resemble those of the more costly high-strength low-alloy steels. Furthermore, spray forming may be an enabling technology for metals such as high silicon steels that are prone to segregation, which leaves the material too brittle to roll into strip. Rapid solidification in spray forming limits segregation and may allow these materials to receive conventional thermomechanical treatments.

Strip Preparation

The INEL spray-forming approach for producing metal strip is depicted schematically in Figure 1. Gas flow through a converging-diverging (de Laval) nozzle creates a localized low-pressure region near the

PRECEDING PAGE BLANK NOT FILMED

14
INTENTIONALLY BLANK



N92 0237

Figure 1. Schematic of INEL approach for spray forming metal strip.

nozzle's throat. Liquid metal is aspirated or pressure-fed into the nozzle through a series of holes, or a slit orifice, that spans the width of the nozzle. The high velocity, high temperature inert gas jet shears and atomizes the metal into fine droplets that are entrained by the gas stream and transported to a rotating drum substrate. In-flight convective cooling followed by conductive and convective cooling at the substrate result in rapid solidification, restricting grain growth and improving product homogeneity by reducing the segregation of impurities that form inclusions. The highly directed two-phase flow can result in the near-net-shape production of strip, as well as complex shapes.

The spray apparatus has been described previously [2]. The design is modular, allowing experimental flexibility for scale-up or the incorporation of specialized components such as a plume diagnostics unit. The apparatus used for continuous strip production consists of a gas manifold and associated electronics for controlling gas flow and temperature, a chamber housing the main spray forming components (induction gas heater, melt tundish, and nozzle), a chamber housing the water-cooled drum substrate, and data acquisition and process control electronics. Process control includes open- or closed-loop computer control of the spray process, laser-based feedback control of strip thickness and surface roughness, and remote video monitoring of the spray process. An in-flight particle diagnostics system is used to simultaneously measure single particle size, velocity, and temperature in the atomized plume. This system measures particle diameters between 5 and 1000 μm using an absolute magnitude of scattered light technique. Velocities of 10 to 100 m/s are measured with a dual beam laser Doppler velocimeter, and particle temperature is measured with a high-speed two-color pyrometry technique.

Bench-scale linear converging/diverging nozzles of our own design were machined in-house from boron nitride. Interchangeable inserts of high purity Al_2O_3 were used in critical areas to minimize erosion of the boron nitride by the molten steel. The throat width, transverse to the direction of flow, was about 25 mm. Mass throughputs were as high as 43 Mg/h per meter of slit width for a slit orifice nozzle operating in the aspiration mode, and 165 Mg/h-m for the same nozzle operating in the pressurized feed mode. A purged argon atmosphere within the spray apparatus minimized slag formation in the melt, surface oxidation of the strip, and in-flight oxidation of the atomized droplets.

The nozzle operated at a static pressure of 206 kPa (30 psia) absolute, measured at the nozzle's inlet. The gas flow field under single-phase flow conditions was mapped using small pitot tube probes. The driving pressure was found to generate supersonic flow conditions; the shock front was in the diverging section near the metal feed location. Gas-to-metal mass ratios typically ranged from 0.1 to 0.5. The gas

and droplet cooled rapidly after exiting the nozzle as the spray plume entrained cool ambient argon. Gas and droplet velocity also decreased after exiting the nozzle, with large droplets responding less to drag effects by virtue of their greater momentum.

The starting material was remelted SAE 1008 hot band. During a typical run, 1.5 kg of steel was induction heated to about 100°C above the liquidus temperature and atomized using argon heated to about 1000°C. The resultant droplets impacted a water-cooled, grit-blasted mild steel drum, producing a strip of metal about 127 mm wide x 3 mm thick.

Spray-Formed Strip Properties

The microstructure of the as-deposited steel was usually fine, equiaxed ferrite with 11 to 45 μm average grain size. The transformation of the microstructure of SAE 1008 steel as it goes from commercial hot band to as-deposited material and finally to hot-rolled product is shown in Figure 2. Note that the average grain size of the as-deposited material is about the same as that of the commercial hot band (~16 μm), but the grains are somewhat more directional, reflecting the heat transfer direction. The grain structure of the spray-formed and hot-rolled material was equiaxed ferrite with ~5 μm grain diameters.

As-deposited density, measured by water displacement using Archimedes' principle, ranged from 88 to 97% of theoretical density with 96% being typical. Full densification of the as-deposited strip was achieved with standard hot deformation processing. Depending upon the sample, hot rolling at 1000 to 1100°C to 30 to 70% thickness reduction was sufficient. Porosity in the as-deposited material was generally highest at the substrate due to high initial quench rates (10^4 to 10^6 °C/s) [3-6]. Thin deposits formed from low density sprays had the highest porosity levels, but also finer grains due to rapid solidification. Low porosities together with fine microstructures were obtained with conditions that favored the formation of dense sprays consisting of small droplets with low solid fractions. The refined and uniform microstructures of thin hot-rolled strip generated under these conditions can be seen in Figure 3. Thick samples (>~9 mm) formed from high enthalpy plumes also had lower porosity levels but coarser as-deposited microstructures. Hot rolling to 35% thickness reduction at 1000°C produced a fine equiaxed ferrite grain structure with an average grain size similar to commercial hot band material. A photomicrograph of strip formed under these conditions, as well as one for commercial hot band, is shown in Figure 4.

As expected, the tensile properties of the spray-formed and hot-rolled low-carbon steel strip reflect the observed grain refinements. Table 1 summarizes the results. The range of values arises from differences

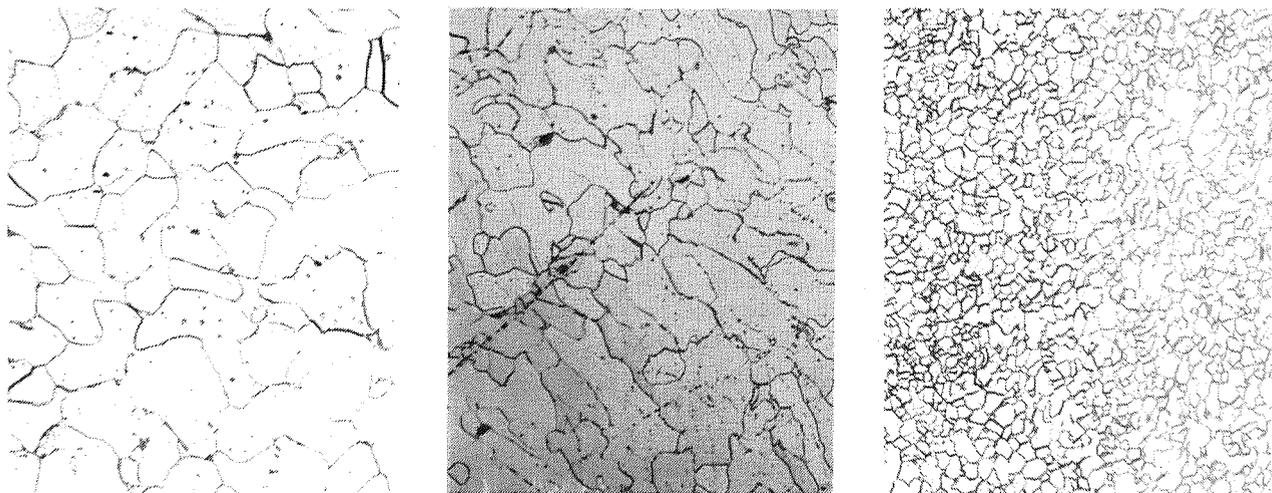


Figure 2. Microstructures of commercial SAE 1008 hot band (left), as-deposited spray-formed strip (center), and hot-rolled spray-formed strip (right). 400X

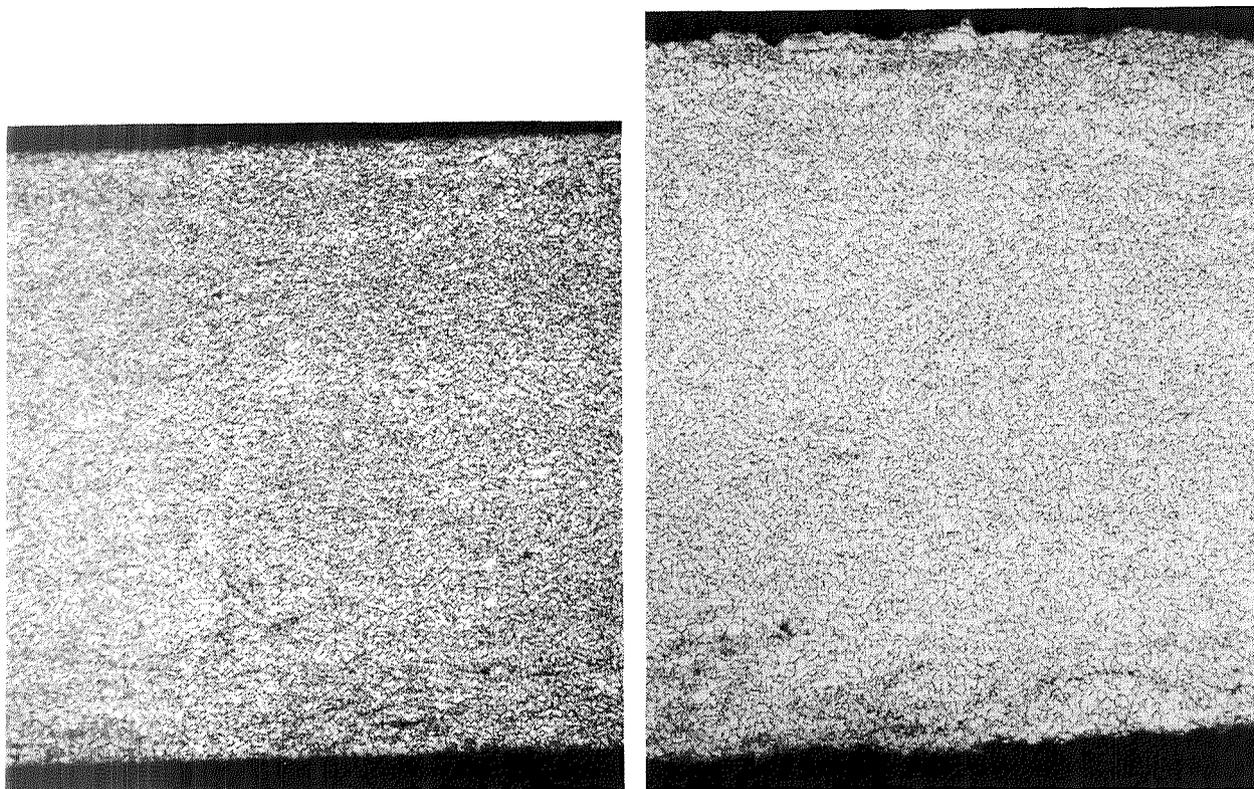


Figure 3. Cross section of thin spray-formed and hot-rolled strip. Left, 100X; right, 35X

in processing parameters (particularly spray conditions). Compared to commercial hot band, yield strengths increased 47 to 64% and ultimate strengths 9 to 63%. The observed reduction in elongation was largely restored by normalizing the samples (heating to 930°C for ~5 min followed by air cooling). Fully annealed samples (heated to 930°C followed by very slow cooling in the furnace) underwent the expected grain growth, with a notable decrease in tensile strength and hardness, and an increase in ductility.

POLYMER MEMBRANES

The unique capabilities of polymer membranes in a wide variety of gas/gas and liquid/liquid separations is well established [7]. The transport properties depend on the membrane's microstructure or "fabric" as well as the physicochemical properties of the polymer and the operating conditions [8]. The microstructure, in turn, is dictated by the fabrication method. The need to fully exploit the potential of existing membrane materials through performance-enhancing fabrication techniques is underscored by the increasingly stringent requirements for air and water purity and waste minimization (see, for example, Title 1 of the Clean Air Act) [9]. To this end, spray forming technology developed for metal coatings was adapted to polymer membrane fabrication. Membranes were spray formed from poly[bis(phenoxy)phosphazene] (PPOP), an inorganic polymer exhibiting exceptional stability in the adverse thermal (>100°C) and chemical (extreme pH) environments frequently encountered in industrial separations [7]. The gas/gas and liquid/liquid separation performance of spray-formed membranes compared favorably with that of similar membranes produced by the conventional method of evaporative knife-casting [10].

Membrane Preparation

Membranes were formed by depositing atomized droplets of linear PPOP dissolved in tetrahydrofuran (THF) onto glass substrates using a bench-scale apparatus developed for spray forming metal. The

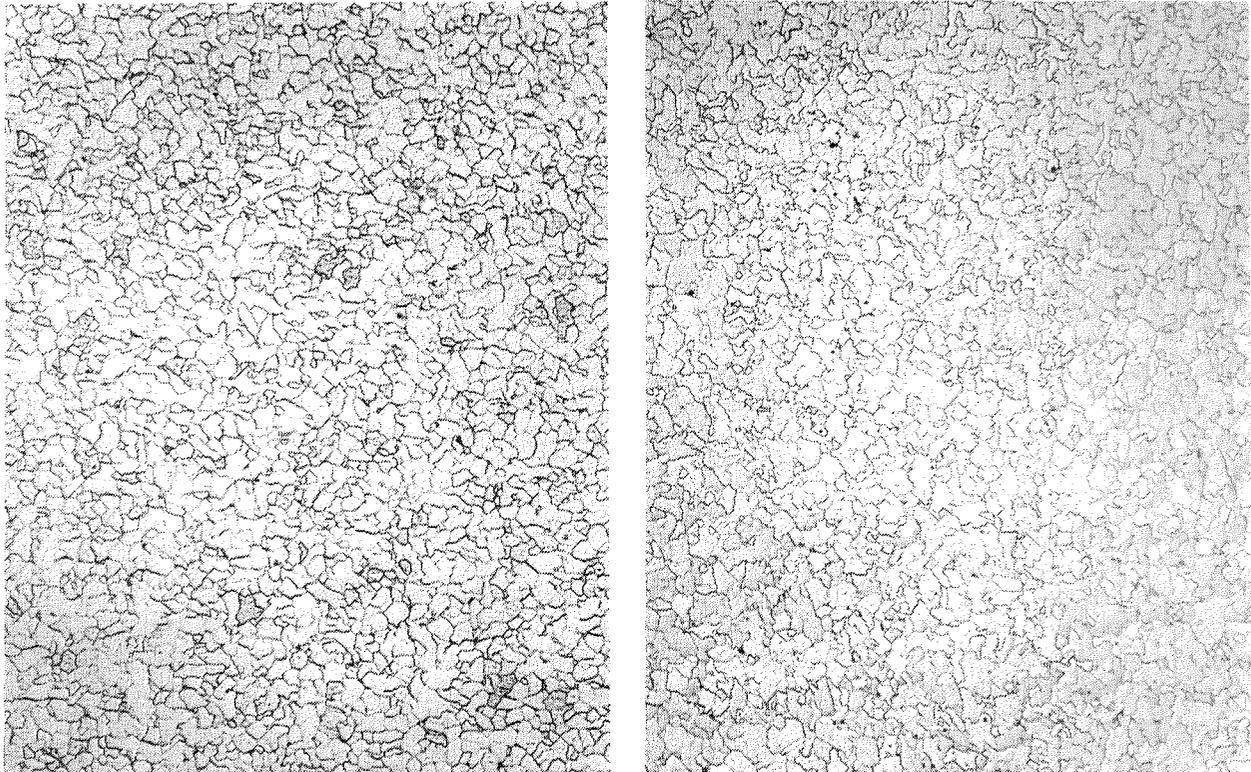


Figure 4. Microstructure of SAE 1008 steel hot band (left) and thick (>~9 mm) spray-formed and hot-rolled strip (right). 100X

Table 1. Tensile Properties of Commercial SAE 1008 Hot Band and INEL Spray Formed and Hot Rolled Strip.

Sample	Yield Strength 0.2% Offset, MPa (ksi)	Ultimate Strength, MPa (ksi)	Elongation, % in 50 mm	Hardness, DPH 100 g Loads
Commercial 1008 Hot Band	197 (28.6)	306 (44.4)	51.8	91
Spray Formed and Hot Rolled	290-324 (42.0-47.0)	334-498 (48.4-72.3)	13.9-37.7	136-160

chemical stability of the polymer allowed the membranes to be sprayed in air; argon was the atomizing gas. The linear converging-diverging (de Laval) nozzle, designed at INEL, was machined from boron nitride.

A 7 wt% solution of linear PPOP in THF was sprayed. The weight average molecular weight of the polymer, measured by gel permeation chromatography, was about 750,000 amu. (3 and 5 wt% solutions having polymer weight average molecular weights exceeding one million amu were also sprayed but gave less satisfactory results.) The solution was warmed to ~45°C to lower its viscosity and poured into the tundish of the nozzle, which operated at a static pressure of 137 kPa (20 psia). The solution was aspirated through six small orifices that spanned the width of the nozzle. Solution throughput was about 0.4 kg/s per meter of nozzle throat width. The corresponding gas-to-polymer-solution mass ratio was about 4. The solution was sheared and atomized, resulting in very fine droplets that were entrained by the gas stream and transported to a moving substrate. Solvent molecules were shed from the atomized droplets during their flight, and the remainder evaporated at the substrate. While control of atomizing gas temperature

could be a convenient means of adjusting the solvent evaporation rate, room temperature argon was used because the equilibrium vapor pressure of THF (145 torr at 20°C) is high enough to allow facile evaporation of the solvent. Upon impact, individual polymer molecules within adjacent droplets interwove while shedding any remaining solvent.

The polymer/solvent spray was deposited onto 8.3 x 8.3 cm glass plates, maintained at room temperature, that were swept through the spray plume at rates yielding membranes 1 to 10 μm thick. A typical 5 μm thick membrane was fully dried and consolidated in about 1 s. The membranes appeared to be coherent and uniform and exhibited good adhesion to the substrate. SEM analysis revealed an asymmetric structure as described below.

The hydrophobic membranes were lifted from the substrate by water immersion, mounted onto a porous test cell support, and edge-sealed using polymer solution. Knife-cast membranes prepared using a standard approach were also floated off their glass substrates into water and mounted onto test cell supports.

Membrane Characteristics and Performance

The microstructure of the spray-formed PPOP membranes was dictated by the interplay of the spray plume and substrate. Experimental parameters such as nozzle geometry and pressure, average molecular weight of the polymer, viscosity and concentration of the polymer solution, evaporation rate of the solvent, and distance to the substrate had the greatest affect on the spray plume characteristics. The speed, temperature, and properties of the substrate also influence the membrane's microstructure.

Scanning electron microscopy (SEM) was used to evaluate membrane surface structure and thickness. Over the width of the glass plates, the membranes appeared homogeneous and of uniform thickness. Close examination revealed that the membranes were asymmetric, with a thin, dense region at the substrate/deposit interface and a relatively thick, uniform build-up of translucent, "spongy" polymer material away from the substrate. Knife-cast membranes, on the other hand, appeared more uniformly dense and transparent.

Gas selectivity was measured using a fully automated mixed gas test cell interfaced to a Hewlett Packard 5190 gas chromatograph [11]. Pervaporation studies were conducted at 65°C using a driving pressure of 200 psig across the membranes. The selectivity of spray-formed and knife-cast PPOP membranes was determined for several acid gas mixtures (10% SO_2 /90% N_2 , 10% H_2S /90% CH_4 , 10% CO_2 /90% CH_4). SO_2 / N_2 mixtures are encountered in industrial exhaust while H_2S / CH_4 represents well gas. The results are given in Table 2. At 80°C, spray-formed membranes had 4 times the selectivity of knife-cast membranes when separating SO_2 from nitrogen; at 130°C the difference increased to about 42 times. Spray-formed membranes had twice the selectivity of similar knife-cast membranes when separating H_2S from methane at 80°C and had 67 times the selectivity at 130°C. Improvements were also observed with spray-formed membranes when separating CO_2 / CH_4 mixtures.

Spray-formed membranes also performed impressively in certain liquid/liquid separations. In pervaporation experiments, spray-formed membranes gave excellent component separation for a mixture of halogenated hydrocarbons and alcohols in water (0.5% methylene chloride, 0.5% chloroform, 0.5% methanol, 0.5% ethanol, 98% water). The alcohol-rich permeate (41.8% ethanol, 58.2% ethanol) contained a trace amount of water. The halogenated hydrocarbon concentration in the permeate was extremely low--below the detection limit in gas chromatographic analysis. Knife-cast membranes gave similar results. However, the flux through spray-formed membranes (2.83 $\text{L}/\text{m}^2\cdot\text{h}$) was appreciably higher than that through knife-cast membranes (0.1-0.4 $\text{L}/\text{m}^2\cdot\text{h}$) of the same thickness tested under the same conditions.

Membrane fabrication via spray forming offers time savings, flexibility, and improved performance over traditional approaches (e.g. knife casting or spin casting). Whereas knife-cast membrane preparation

Table 2. Component Selectivity Data for Spray-Formed and Knife-Cast PPOP Membranes.

Gas Mixture	Temperature (°C)	PPOP Selectivity	
		Spray-Formed	Knife-Cast
10% SO ₂ /90%N ₂	80	71:1	18:1
10% H ₂ S/90% CH ₄	80	15:1	7:1
10% CO ₂ /90% CH ₄	80	4.5:1	3.5:1
10% SO ₂ /90% N ₂	130	344:1	7.2:1
10% H ₂ S/90% CH ₄	130	303:1	5:1

required hours, spray-formed membranes were prepared in seconds. The flexibility gained by spray forming membranes to near-net shape not only greatly reduces production costs by eliminating unit operations, but also allows membranes with complex shapes, which are difficult or impossible to manufacture by conventional approaches, to be produced in a straightforward manner. The ability to tailor membrane microstructure to specific separation processes by varying the spray and substrate parameters mentioned previously enhances performance.

SPRAY-FORMED TOOLING

The recent explosion of interest in rapid prototyping technology is fueled in part by the restructuring of today's marketplace. Successful competition in global markets will require the ability to carry a design concept through the prototype stage to the production stage faster and at lower cost than ever before. The ability to generate plastic and wax models of prototype parts with high dimensional accuracy via selective laser sintering [12], stereolithography [13], ballistic particle manufacturing [14], and other approaches is now a reality. However, it is generally accepted that the rapid production of prototype parts from engineered materials--materials that will actually see service--is the prime long term goal [15]. Methodologies that can rapidly produce specialized tooling, such as molds and dies, would satisfy this goal when used with conventional approaches such as injection molding, compression molding, and die casting.

Presently, complex molds, dies, and related tooling are expensive and time consuming to make. They can easily exceed \$200K in cost and require months to fabricate. Researchers at the INEL have recently begun to develop spray forming technology to produce specialized tooling by spray depositing molten metal droplets onto patterns made from plastics, waxes, clay, or other easy-to-form materials. This approach could provide a unique opportunity for simplifying production of complex tooling, thereby substantially reducing its cost. Rapid solidification enables patterns made from plastics, waxes, clays, etc. to be used despite their low softening temperatures, while near-net-shape capability allows objects with complex shapes to be made easily. The semispherical shell shown in Figure 5 was generated by spray depositing molten tin directly onto an inflated party balloon.

Experimental

To form a mold, die, etc., liquid metal was pressure fed into a venturi-like nozzle transporting high velocity (mach number ~ 1.5) argon at temperatures above the liquidus temperature of the metal to be sprayed. Kinetic energy transfer from the gas overcame the relatively strong surface tension forces of the liquid metal, resulting in finely atomized metal droplets. The droplets were entrained in a directed two-phase flow, quenched, and deposited onto a moving plastic pattern having the desired shape and surface texture. The main spray-forming components (spray nozzle, liquid metal reservoir, gas heater, and pattern)

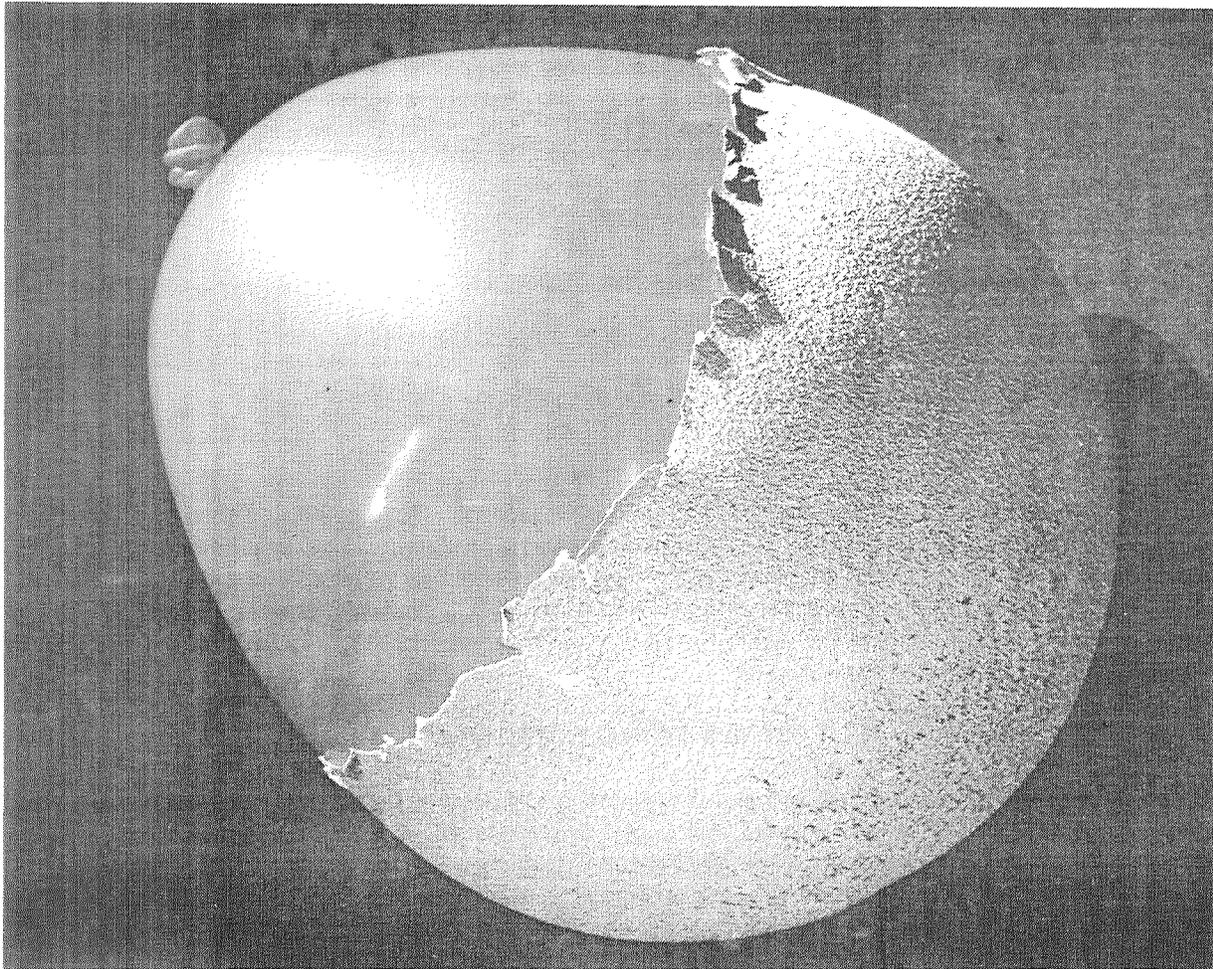


Figure 5. Party balloon spray-coated with tin emphasizes rapid solidification and near-net-shape capability of spray forming.

were housed in an argon-purged chamber to limit the detrimental effects of oxide formation. All spray components were designed and constructed in-house.

The nozzle/metal feed assembly was designed to produce sprays of relatively fine droplets having a narrow size distribution. These conditions offer greater flexibility for controlling droplet temperature, momentum, and flow pattern, as well as deposit microstructure. Bench-scale nozzles having transverse throat widths of 17 mm were typically operated at gas-to-metal mass ratios (for tin) of about 10 with metal throughputs of about 0.5 kg/s per meter of nozzle throat width.

Single-phase gas flow field diagnostics were used to map the static pressure and gas velocity profiles within the nozzle's flow channel. Size analysis of solidified droplets was conducted using standard wet and dry sieving methods.

A quasi one-dimensional computer model was used with the diagnostics results to guide component optimization as well as development of algorithms for process control. The model simulated the entire nozzle and free jet (plume) regions with full aerodynamic and energetic coupling between the metal droplets and the transport gas, and with coupled liquid injection into the gas stream.

Results and Discussion

The ultimate goal is fabrication of complex tooling from tool forming and hardfacing stainless steel alloys and composite materials. At the time of this publication, however, development of spray-forming tooling at the INEL is in its early stages. Several low-melting point alloys of zinc and tin have been tested with very encouraging results. An example is given in Figure 6, which shows a metal mold produced in about 5 min by spray-forming molten tin on a plastic (low-density polyethylene) pattern having a butterfly shape. The pattern was not damaged despite the fact that the temperature of the molten metal within the crucible (300°C) greatly exceeded the melting point of the pattern (~100°C). Replication of surface features, including fine scratches in the pattern, was excellent. The surface of the mold was mirror-like and free of solidification shrinkage defects, indicating that replication of the pattern's surface texture also was very good. Patterns of a variety of other plastics, including acrylic, polycarbonate, and polystyrene, have also given good results.

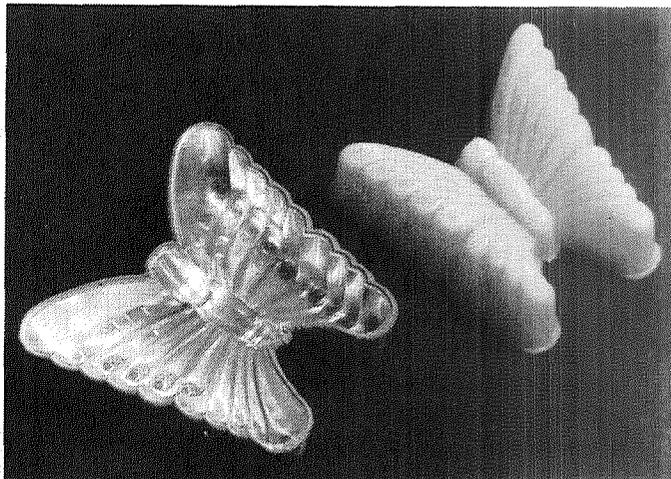


Figure 6. Metal mold shell (left) was produced in about 5 min by spray depositing tin on a plastic pattern (right).

A photomicrograph of a sectioned mold, given in Figure 7, illustrates the refined grain structures that can be obtained using this rapid solidification process. The as-deposited grain structure was equiaxed with a fairly narrow range of fine (~6-16 μm) grain sizes--much finer than the massive grains found in cast objects. As-deposited density, measured by water displacement using Archimedes' principle, was typically 88 to 95% of theoretical.

The molten metal used to produce the deposit was very finely atomized. Unconsolidated powder was collected and analyzed by wet and dry sieving through fine mesh screens. The mass median diameter, volume mean diameter, and Sauter mean diameter of the powder were, respectively, 23 μm , 31.3 μm , and 23.2 μm . The geometric standard deviation ($\sigma_g = (d_{84}/d_{16})^{1/2}$) as determined from a log-normal plot) was 1.5, indicating a narrow droplet size distribution in the spray plume. SEM analysis revealed that nearly all the particles were spherical.

An important advantage of spray forming molds, dies, etc. is the ability to use patterns made from easy-to-shape materials such as plastic, wax, or clay, even though the softening point of these materials may be well below the crucible temperature of the molten metal. Since plastic and wax prototype models can now be produced using CAD-based systems, spray forming could develop into a complementary approach for generating specialized tooling for manufacturing prototype parts from engineered materials. The reduced time and cost of these molds/dies would allow rapid design verification and enable new designs and technology to enter the marketplace more quickly.

ACKNOWLEDGMENTS

We gratefully acknowledge significant contributions of Ray Berry in modelling multiphase flow behavior, heat transfer, and solidification phenomena; Denis Clark in substrate development, process control, and modelling efforts in spray deposition; James Fincke and David Swank in particle and gas flow field diagnostics; and Eric Peterson in membrane characterization and testing. This work was supported by the U.S. Department of Energy, Office of Conservation and Renewable Energy, Office of Industrial Technology, and by the EG&G Idaho Laboratory Directed Research & Development Program under DOE Idaho Field Office Contract DE-AC07-76ID01570.

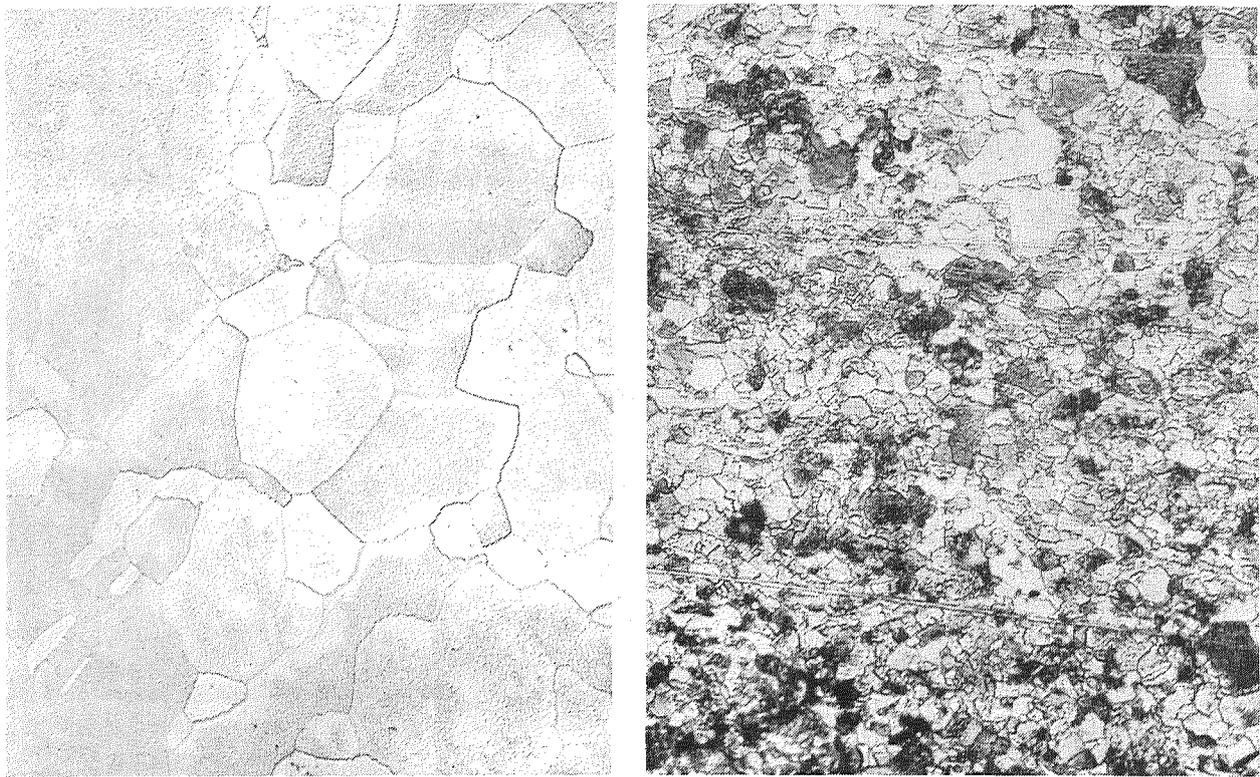


Figure 7. Microstructures of cast tin (left) and spray-formed tin mold of Figure 6 (right). 400X

REFERENCES

1. *Report of the National Critical Technologies Panel*, William D. Phillips, Chairman, The National Critical Technologies Panel, Arlington VA, March 1991.
2. J. F. Key, R. A. Berry, D. E. Clark, J. R. Fincke, and K. M. McHugh, *Development of a Spray-Forming Process for Steel. Final Program Report*, Dec. 1991 (Grant No. DE-AC07-76ID01570).
3. A. R. E. Singer, *Powder Metal.* **25** (4), 195 (1982).
4. P. Predecki, A. W. Mullendore, and N. J. Grant, *Trans. Metall. Soc. AIME* **233**, 1581 (1965).
5. P. Duwez and R. H. Willens, *Trans. Metall. Soc. AIME* **227**, 362 (1963).
6. R. C. Ruhl, *Mater. Sci. Eng.* **1**, 313 (1967).
7. S. A. Leeper, et al. *Membrane Technology and Applications: An Assessment*, EGG-2282, U.S. DOE Contract No. DE-AC07-76ID01570, Feb. 1984.
8. R. R. McCaffrey and D. G. Cummings, *Separation Science and Technology* **23** (12,13), 1627 (1988).
9. *Spectroscopy* **6** (9), 24 (1991).
10. K. M. McHugh and E. S. Peterson, to be published.
11. E. S. Peterson, M. L. Stone, R. R. McCaffrey and D. G. Cummings, *Separation Science and Technology* (in press).
12. C. R. Deckard, *Manufacturing Processes, Systems, and Machines: 14th Conference on Production Research and Technology*, S. K. Samanta, Ed., NSF, Ann Arbor, MI, 1987.
13. L. E. Weiss, E. L. Gursoz, F. B. Prinz, P. S. Fussell, S. Mahalingam, and E. P. Partick, *Manufacturing Review* **3** (1), 40 (1990).
14. D. Hauber, *Manufacturing Processes, Systems, and Machines: 14th Conference on Production Research and Technology*, S. K. Samanta, Ed., NSF, Ann Arbor, MI, 1987.
15. S. Ashley, *Mechanical Engineering*, **113** (4), 34 (1991).

Film Fabrication Technologies at NREL

Robert D. McConnell
National Renewable Energy Laboratory
Golden, CO 80401

S2-27
150472
N 93-25563

ABSTRACT

The National Renewable Energy Laboratory (NREL) has extensive capabilities for fabricating a variety of high-technology films. Much of the in-house work in NREL's large photovoltaics (PV) program involves the fabrication of multiple thin-film semiconducting layers constituting a thin-film PV device. NREL's smaller program in superconductivity focuses on the fabrication of superconducting films on long, flexible-tape substrates. This paper focuses on four of NREL's in-house research groups and their film fabrication techniques, developed for a variety of elements, alloys, and compounds to be deposited on a variety of substrates. As is the case for many national laboratories, NREL's technology transfer efforts are focusing on Cooperative Research and Development Agreements (CRADAs) between NREL researchers and private industry researchers. The reader is encouraged to consider and explore the application of these film-fabrication technologies for their own needs by contacting the author or the principal scientists at NREL, referenced below.

INTRODUCTION

Chemical deposition and physical deposition are the two principal classifications for a multitude of film-fabrication technologies, whether the films are thin (about 1 micron thick or less) or not (thicker than 1 micron). Among the chemical deposition technologies used by NREL researchers are several chemical vapor deposition techniques, liquid phase epitaxy, solution growth, and several electrochemical techniques. Among the physical deposition technologies are several sputtering techniques, physical vapor deposition using electron beams or resistance heating, molecular beam epitaxy, and melt coating. These techniques or their variations have been used at NREL to deposit a variety of metals, compounds, and alloys as amorphous, polycrystalline, or crystalline films. Among the critical factors in identifying an appropriate deposition technology are the control precision needed for the amounts of the elements used in the films—especially important in the case of semiconducting compounds and their dopants—the effect of the deposition technology on the crystallinity of the final films, and the deposition rate—an especially important parameter for production costing.

A historical description of NREL's film fabrication technologies and the materials deposited is in reference 1. The deposition techniques at NREL include electron beam evaporation, physical vapor deposition, magnetron (DC and RF) sputtering, ion-beam sputtering, electrodeposition of multielement compounds, molecular beam epitaxy, liquid-phase epitaxy, electrolytic plating, electroless deposition, plasma-enhanced chemical vapor deposition, photochemical vapor deposition, hot-wire catalytic deposition, melt sheet growth, and melt coating. Several of these techniques, or their variations, have been used to deposit films of Si, amorphous Si, amorphous Ge, amorphous C, Au, Sn, Zn, Ag, Mo, W, GaAs, GaAlAs, GaInP₂, GaAsP, InP, GaP, CuInSe₂, CdS, In₂O₃(Sn), SnO₂(Sb), ZnO(Al), Ta₂O₅, TiO₂, MgF₂, ZnS, MoO₃, WO₃, diamond-like carbon films, YBaCuO compounds, BiSrCaCuO compounds, TlBaCaCuO compounds, other refractory metals, other metal oxides, silicon alloys, and amorphous-silicon alloys.

NREL also has a long history of substantial funding (50% of NREL PV funds) for industrial research in photovoltaics; and much of that research is devoted to manufacturing processes, including film fabrication research for thin-film devices. Subcontract research to industry is a legitimate technology transfer activity at NREL, although industry itself is supported to develop technology, and no "transfer" takes place.

This paper will highlight four of NREL's film fabrication technologies and their achievements within NREL's PV and superconductivity programs. Three are chemical deposition technologies that have proven to be advantageous for the controlled deposition of multielement compounds and alloys. The fourth is a physical deposition technology with known advantages for use in production. NREL's in-house research capabilities are available to US PV companies, superconductivity companies, and other companies interested in film-fabrication technologies.

METALORGANIC CHEMICAL VAPOR DEPOSITION

J. Olson and coworkers successfully fabricated record-high-efficiency photovoltaic devices using a custom-built metalorganic chemical-vapor-deposition system (2, 3). They grew their epitaxial layers of GaInP₂ and GaAs in an atmospheric-pressure MOCVD vertical reactor, and achieved efficiencies for the conversion of sunlight to electricity as high as 29% (active area), a record for this technology. The system was built with a standard run-vent configuration. The Ga, In, P, As, and Zn sources were compounds of trimethylgallium, trimethylindium, phosphine, arsine, and diethylzinc. Carbon doping came from a CCl₄ source consisting of reagent-grade CCl₄ housed in a conventional bubbler, with Pd-diffused H₂ used as the carrier gas. Growth was done on zinc-doped GaAs single-crystal substrates.

During their work on improving device performance, this group developed an in-situ growth-rate measurement technique using diffuse reflectance normal to the growth surface. The technique allows for measurement of real-time surface roughness, in addition to growth rate. The period of interference fringes obtained from the diffuse reflectance was compared with cross sections measured by a scanning-electron microscope and shown to be inversely proportional to the growth rate.

HOT-WIRE CHEMICAL VAPOR DEPOSITION

R. S. Crandall and coworkers developed a different chemical vapor deposition technique for amorphous silicon (4). Called hot-wire-assisted chemical vapor deposition (HWCVD), it appears to be capable of producing hydrogenated amorphous silicon with properties superior to those fabricated by the more conventional plasma-enhanced chemical vapor deposition (PECVD). The HWCVD films are deposited using a hot filament temperature of approximately 1,900°C, to dissociate the silane gas, and a chamber pressure low enough to minimize gas collisions before the film precursors hit the substrate. The amorphous silicon films are fabricated using a SiH₄ gas flow rate of 20 sccm and either glass or silicon substrates held at varying temperatures between 40°C and 450°C. Deposition rates tended to be as much as four times faster than PECVD rates; rapid deposition rates are critical to the cost effectiveness of production technologies. Figure 1 is a schematic of this process.

ELECTRODEPOSITION

Instead of chemical vapors, R. Noufi, R. Bhattacharya, and coworkers developed an electrodeposition technique using chemical solutions for depositing superconducting films (5). Electrodeposition is an electrochemical process where materials (e.g., metals or oxides) in the proper stoichiometry are deposited on conducting substrates from chemical solutions containing the ions of interest (e.g., Cu²⁺, Y³⁺, Ba²⁺, in the case of one of the high-temperature superconductors). As is almost always the case for electrodeposition, one of the two phases contributing to an interface of interest will be a liquid electrolyte, which is merely a phase through which charge is carried by the movement of ions. The second phase at the boundary is a solid electrode (the substrate in this case), which is the phase through which charge is carried by electron movement. In general, when the potential of an electrode is moved from its equilibrium value toward negative potentials, the substance that will be reduced first ($\text{Cu}^{2+} + 2e^- = \text{Cu}^0$) is the one with the least negative redox potential. For example, in a solution containing Cu²⁺, Y³⁺, and Ba²⁺, the Cu²⁺ is reduced first, followed by the reduction of Y³⁺ as the potential of the electrode is made more negative. All three ions can be deposited on the surface of the electrode when the potential is negative enough. The relative concentrations of the constituents in the deposited films are empirically determined by the concentrations of the ions in the solution.

There are other families of superconducting compounds involving Ba, Sr, Ca, Cu, O and Tl, Ba, Ca, Cu, O (6). NREL researchers have successfully electrodeposited these compounds in their appropriate stoichiometries. A significantly high rate of fabrication, of the order of microns per minute, makes electrodeposition a likely candidate for large-scale, cost-effective manufacturing of superconducting wires or tapes. Figure 2 is a schematic of this process for the fabrication of superconducting tapes made from thallium compounds.

PILOT LINE SPUTTERING

Many PV researchers have concluded that sputtering does too much damage to films during fabrication to yield the optimum electronic properties needed for high-efficiency photovoltaic devices. Yet T. Coutts and coworkers adapted this deposition technology, with known production advantages to a high-quality solar-cell design destined for space applications (7). Using Ar sputtering gas, InP solar cells were fabricated using two sputter guns and depositing successive layers of indium-tin-oxide (ITO). The first thin layer provided a shallow homojunction in the InP substrate, while the second provides part of an antireflection coating needed for high performance of the completed solar cell. Additional layers, patterned electrical contacts, and a post-deposition heat treatment provide controlled changes in the electrical and optical properties of the solar cell and complete the processing of the final cell. Thirty-two solar cells, each 4 cm² in area, were completed in the pilot production. The histogram of efficiencies ranged from about 15% to over 16%. The project demonstrated that the sputtering process can be used for these relatively-large-area devices and that the highest efficiency, 16.2%, is comparable to the highest reported from another production method. Further, the process can be configured for in-line production rather than batch production, characteristic of many CVD production processes.

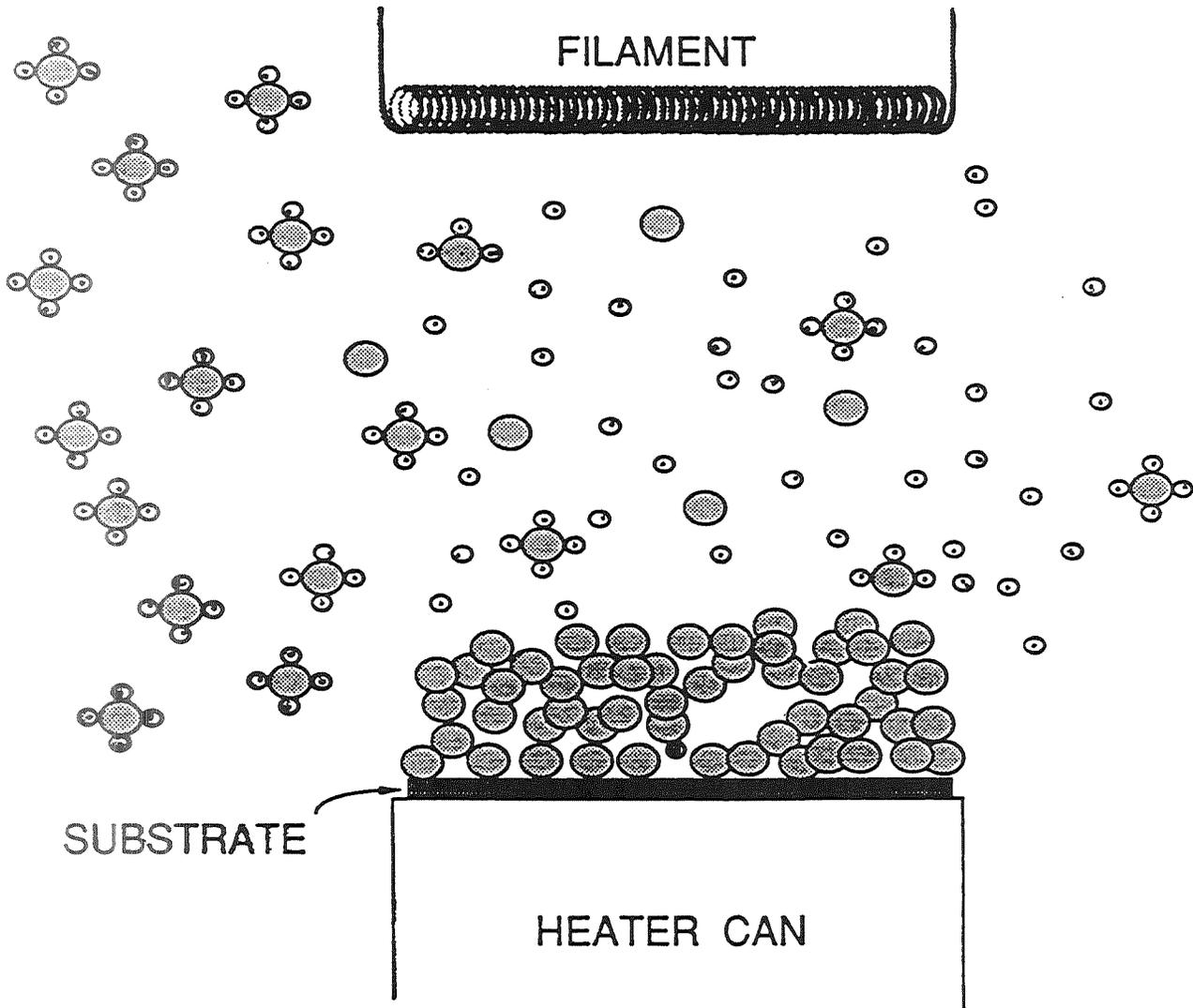
CONCLUSION

This paper has briefly described four film fabrication technologies developed by NREL researchers for either multilayer PV devices or high temperature superconducting coatings. Each of the research groups mentioned above has achieved distinction within its respective research community through the development of these particular film-fabrication processes for high-efficiency PV devices or superconducting coatings. For U.S. PV companies, superconductivity companies, or other companies interested in high-technology films and coatings, these and other NREL research groups are resources available through formal cooperative R&D agreements, such as CRADAs.

REFERENCES

- (1) SERI: *Your Laboratory for Materials Research, Processing, and Characterization*, SERI/SP-21-3600, (1990)
- (2) J. M. Olson, S. R. Kurtz, A. E. Kibbler, and P. Faine, "Recent Advances in High Efficiency GaInP₂/GaAs Tandem Solar Cells," Proceedings of the 21st IEEE Photovoltaics Specialists Conference, (1990)
- (3) A. E. Kibbler, S. R. Kurtz, and J. M. Olson, "Carbon Doping and Etching of MOCVD-grown GaAs, InP, and Related Ternaries using CCl₄," *Journal of Crystal Growth*, **109**, p. 258, (1991)
- (4) A. H. Mahan, J. Carapella, B. P. Nelson, I. Balberg, and R. S. Crandall, "Deposition of Device Quality, Low H Content Amorphous Silicon," *J. Appl. Phys.* **69** (9), P. 6728, (1991)
- (5) R. N. Bhattacharya, R. Noufi, L. L. Roybal, R. K. Ahrenkiel, P. Parill, A. Mason, and D. Albin, *Science and Technology of Thin Film Superconductors 2*, Editors R. D. McConnell and R. Noufi, Plenum Press, New York, P. 243, (1990)
- (6) R. N. Bhattacharya, P. A. Parilla, R. Noufi, P. Arendt, and N. Elliott, "YBaCuO and TlBaCaCuO Superconductor Thin Films via an Electrodeposition Process," *Journal of the Electrochemical Society*, **139**, No. 1, p. 67, (1992)
- (7) T. A. Gessert, X. Li, T. J. Coutts, and N. Tzafaras, "Pilot Production of 4 cm² ITO/InP Photovoltaic Cells," Proceedings of the 21st IEEE Photovoltaic Specialists Conference, (1990)

FIGURE 1. HOT-WIRE CHEMICAL VAPOR DEPOSITION

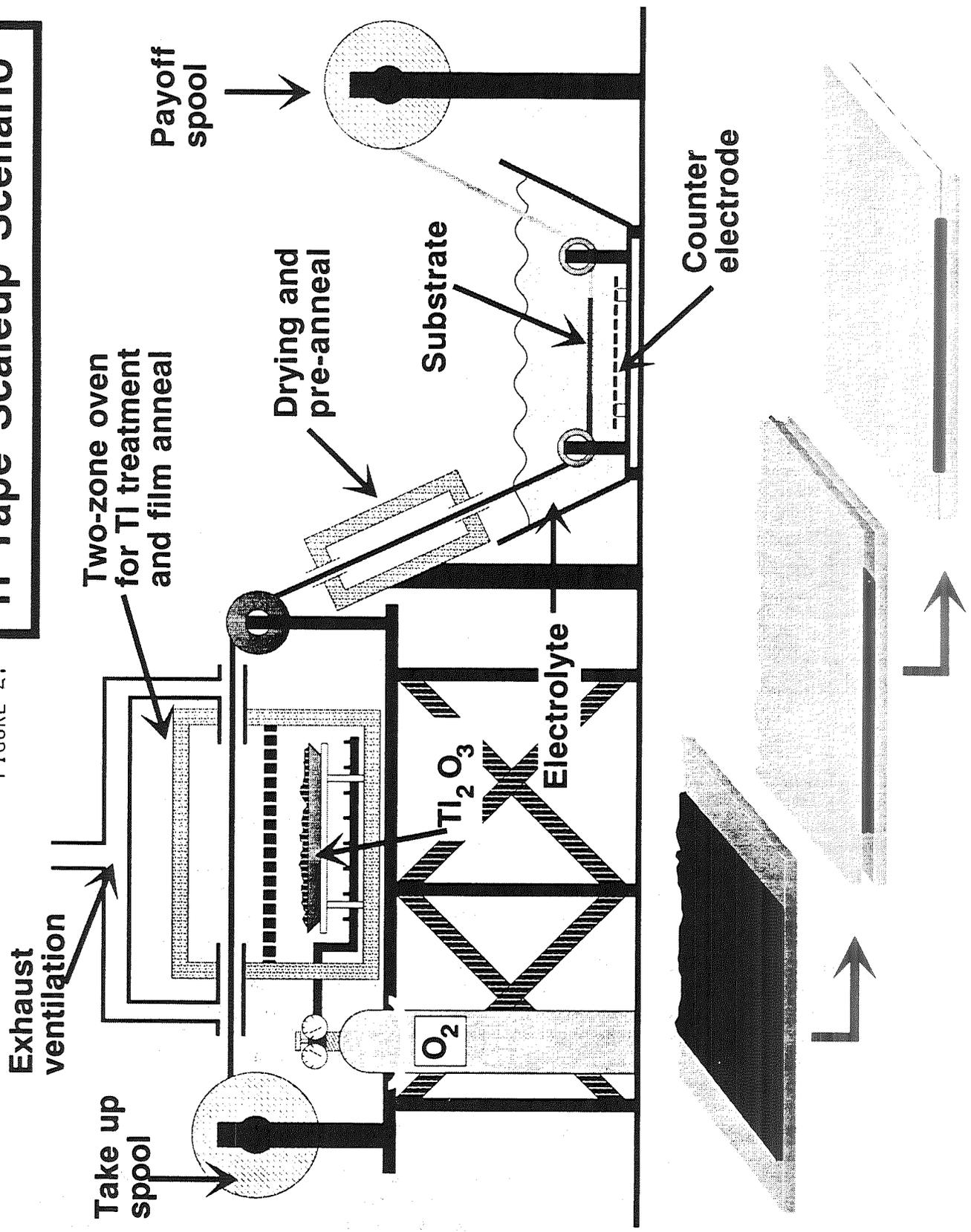


	SILANE (SiH₄)		SiH₃
	HYDROGEN		SILICON

Deposition Rate $\approx 10 \text{ \AA/second}$

TI Tape Scaleup Scenario

FIGURE 2.



S3-27
150473

l-11

**OPTICAL and SCRATCH RESISTANT PROPERTIES of
DIAMONDLIKE CARBON FILMS DEPOSITED
WITH SINGLE and DUAL ION BEAMS**

N 9 3 - 2 5 5 6 4

Michael T. Kussmaul
Sverdrup Technology, Inc.
Brook Park, Ohio

Michael S. Bogdanski*¹
Case Western Reserve University
Cleveland, Ohio

Bruce A. Banks
Michael J. Mirtich
NASA Lewis Research Center
Cleveland, Ohio

ABSTRACT

Amorphous diamondlike carbon (DLC) films were deposited using both single and dual ion beam techniques utilizing filament and hollow cathode ion sources. Continuous DLC films up to 3000 Å thick were deposited on fused quartz plates. Ion beam process parameters were varied in an effort to create hard, clear films. Total DLC film absorption over visible wavelengths was obtained using a Perkin-Elmer spectrophotometer. An ellipsometer, with an Ar-He laser (wavelength 6328 Å) was used to determine index of refraction for the DLC films. Scratch resistance, frictional and adherence properties were determined for select films. Applications for these films range from military to the ophthalmic industries.

INTRODUCTION

Extensive resources by a large number of laboratories have been dedicated to diamond film research. The attraction for diamond films is easy to understand. Diamond is the hardest known material, an excellent insulator, and has a high thermal conductivity. However, high substrate temperatures ($\approx 1000^{\circ}\text{C}$) are required for the formation of these films. Thus, their usefulness is restricted to specialized applications and to relatively small surface areas. Hydrogenated amorphous DLC films (1) however, can be deposited at low temperatures, thus attracting numerous applications which are unapproachable by the high temperature diamond film deposition technique. DLC films can be deposited at room temperature using several plasma generating techniques as detailed by Robertson (2). With commercially available ion source systems, these deposition techniques are easily configured to coat large and/or unique samples. Ideally the only constraint regarding film deposition using ion sources is the size and pumping limitations of the vacuum facilities.

NASA Lewis Research Center is performing spinoff technology development to improve the characteristics of DLC films made by dual ion beam deposition for the purposes of coating plastic ophthalmic lenses. High visible wavelength transmittance, scratch-resistance, film hardness, and acceptable index of refraction are required for such applications.

Diamondlike Carbon Film Deposition Apparatus and Procedure

Two different dual ion beam systems were used to deposit DLC films. The first system was used by Mirtich (1). It consisted of a 30 cm ion source with its grids masked down to 10 cm in diameter and an 8 cm ion source (figure 1). Argon flowed through a hollow cathode and a neutralizer to create a plasma dis-

*NASA Resident Research Associate at Lewis Research Center.

charge. Methane was introduced into the discharge chamber to provide a source of carbon for DLC film deposition. A second system used to deposit DLC consisted of a 15 cm diameter source and 8 cm diameter source. The 15 cm cathode filament ion source with its extraction grids masked to 10 cm diameter, is used to directly deposit DLC films using methane.

Prior to film deposition, samples are cleaned with soap solution in an ultrasonic bath, rinsed with deionized water and dried with nitrogen. The specimens to be coated are then placed in vacuum and cleaned using xenon ions at an energy of 500 eV for approximately two minutes. The total beam energy (the sum of the discharge voltage and the screen grid voltage) is kept at approximately 125 eV. Current densities at these conditions are approximately $60 \mu\text{a}/\text{cm}^2$ in the vicinity of the sample. The deposition rate of the DLC films on fused quartz is approximately $71 \text{ \AA}/\text{min}$.

Robertson and O'Reilly (3) have shown that a mixture of sp^2 (trigonal bonding associated with graphite) and sp^3 (tetrahedral bonding characteristic of diamond) form into graphitic clusters which are then bonded into a larger sp^3 matrix. It was also shown that the sp^2 clusters control the electronic properties in the DLC film while the sp^3 bonding is responsible for its mechanical properties. Angus and Wang (4) also discuss the role of atomic hydrogen content in DLC films as a method of increasing the sp^3 coordination while reducing the sp^2 bonding.

Mirtich, et al (5), showed that a dual ion beam system with energetic argon ions in the second ion source could produce clearer DLC films than possible with a single ion direct deposition source. During dual-beam depositions an 8 cm cathode filament ion source with its extraction grids masked to 1 cm, is used to direct a beam of hydrogen, argon or xenon ions at the substrate. The current density due to hydrogen ions ranged from $80\text{-}300 \mu\text{a}/\text{cm}^2$.

Optical Properties of Films

The spectral transmittance and reflectance of the DLC films deposited on fused quartz were obtained using a Perkin-Elmer lambda 9 spectrophotometer. The spectral absorptance is calculated based on the measured transmittance and reflectance. Figure 2 shows the transmittance for DLC films deposited with both the single and dual ion beam systems. The transmittance of the DLC film deposited with the dual ion beam system is greater for all wavelengths compared to the single beam film. A 500 \AA thick film also is 90% transmitting at wavelengths greater than 7000 \AA .

Figure 3 shows the results of various techniques which were tested in an effort to increase the DLC film transmittance at a wavelength of 5000 \AA . The values represent the transmittance at 5000 \AA (which approximates the peak sensitivity of the human eye (6)) for DLC films on fused quartz substrate for various dual beam gaseous deposition conditions. The 30 cm source for conditions 1 through 4 used a hollow cathode to produce a plasma of argon and methane for DLC film deposition. If the hollow cathode is replaced by a filament cathode the need to use Argon in the operation of the ion source during the DLC deposition process is eliminated. The improvement in DLC film transmittance as a result of this alteration is shown in condition 5 (Figure 3). Clearly a higher transmittance results when the dual beam system is used, and especially when hydrogen gas is used in the second source. The use of a pure hydrogen ion beam in the 8 cm ion source produced a DLC film with a transmittance of 84%.

Index of refraction measurements were made with an Ellipsometer II system manufactured by Applied Materials, Inc.. This system uses a 2 milliwatt helium-neon laser, with a wavelength of 6328 \AA and at a 70° angle of incidence with the surface. This ellipsometer produces data which is used to calculate both the index of refraction and the film thickness. The calculated film thickness is compared to the value obtained using a Dektak surface profilometer to determine the reliability of the calculated index of refraction. The techniques showed agreement within approximately 70 \AA .

A DLC film's index of refraction is important because this property determines the required film thickness for use in an anti-reflective stack coating. As the difference between two materials' indices of refraction increases, less of each material is required to produce an anti-reflective coatings. This smaller amount of material (thinner film) would correspond to a higher transmittance because of less light attenuation and to a reduction in film deposition time. both of these factors (less material, less time) translate to a lower cost to produce anti-reflective coatings.

The DLC films made with the direct deposition technique using methane in a single ion source produce films which have an index of refraction of 2.0. When a hydrogen beam is directed at the sample during film deposition the index of refraction decreases to values of 1.75 to 1.8 as the current density of the hydrogen beam in the second ion source is increased.

Scratch Resistance and Adherence

Diamond stylus scratch tests were performed on both a DLC coated and an uncoated fused quartz plate. The DLC coating had a thickness of 1650\AA , and was deposited using a single ion beam source. The plates were ultrasonically cleaned in acetone, rinsed with pure ethanol followed by deionized water, then dried in air. A diamond stylus with a hemispherical radius of $841\ \mu\text{m}$ was slid along the plate at a rate of 10 mm/min. A progressive normal load of 0 to 25N was applied at a rate of 100 N/min. These parameters generated a total scratch length of 2.5 mm with an effective load-displacement relationship of 10 N/mm. The frictional force and acoustic emission were monitored during the scratch tests. The acoustic emission signal is an accurate indicator as to when fracture initiates for brittle materials, such as fused quartz.

Figure 4 summarizes the results from the scratch tests on the coated and uncoated plates. The acoustic emission signal indicates that fracture initiated on the uncoated quartz plate began to fracture at approximately 13 N normal load whereas the DLC coated plate did not begin to fracture until approximately 21 N normal load. The frictional force was relatively linear in both cases, but was notably lower for the coated plate. The average friction coefficient (frictional force divided by normal force) for the test duration was 0.04 ± 0.01 for the uncoated plate and 0.03 ± 0.01 for the DLC coated plate.

This reduction in friction is the most likely reason that the DLC coated plate withstood a higher normal load before fracturing than the uncoated plate. A spherical stylus sliding along a flat plate generates a maximum tensile stress at the trailing edge of the sliding contact (7). It has been shown that this trailing edge tensile stress is the cause of failure in brittle materials (ref 8). This tensile stress is amplified by the friction coefficient between the stylus and the plate. Therefore a reduction in friction coefficient also reduces the tensile stress at the trailing edge of the contact.

Abrasion tests were performed on single ion beam DLC films deposited on fused quartz by rubbing SiO_2 particles ($\approx 80\ \mu\text{m}$ particle size) over the surface. These particles were rubbed on the surface by hand. Figure 5 shows how the DLC film coated side was protected from the abrasive particles while scratches can be seen on the uncoated portion of the sample.

Adherence of DLC films deposited on fused quartz using both the single and dual beam methods was measured by Mirtich (9). The adherence of these films were as good as the maximum adherence of the Sebastian Adherence Tester used to make the measurement ($5.5 \times 10^7\ \text{N/m}^2$ or 8000 psi), regardless of the method of deposition. The film adherence was so great in fact that often portions of the quartz lifted off leaving the DLC film intact.

Potential DLC Applications

A variety of companies in the United States were surveyed in 1989 by the regarding the potential commercial applications of DLC films. Table I lists the 12 applications which were most highly rated of the 39 questionnaire responses received. In addition, Table II lists of the focussed applications which are being explored at NASA Lewis Research Center for which DLC films would be well suited.

SUMMARY

DLC films are hard, very adherent and can be deposited at room temperature. These films may prove to be beneficial in technological areas seeking to improve the scratch resistance of materials through the use of novel thin coatings. The availability of ion beam systems (whether dual or single beam) facilitates configuring systems for deposition onto large, non-standard shapes. Single beam DLC films are already being used as scratch resistant coatings on sunglass lenses which are available commercially since high transparency is not as critical. The addition of a second ion source has been shown to improve film transmittance, thus increasing the usefulness of these films to other eyewear and optical surfaces (automotive windows, scanners, etc.). The index of refraction for DLC films is reduced when a dual ion system

employing a hydrogen beam is used. A DLC coated quartz plate has shown a superior ability to resist fracture under sliding in comparison to an uncoated quartz plate. Thus, the use of a fairly thin DLC film has the potential to provide improved scratch resistance beyond the capability of the substrate onto which it is deposited.

REFERENCES

1. Mirtich, M.J., et al, Ion Beam and Plasma Methods of Producing Diamond-Like Carbon Films, Thin Solid Films, 131, pp. 245-254, (1985).
2. J. Robertson, Properties of Diamond-Like Carbon, Surface Coatings Technology 00, 1, (1992).
3. J. Robertson and E.P. O'Reilly, Phys. Rev. B 35, 2946 (1987).
4. J.C. Angus and Y. Wang, Diamond and Diamond-Like Films and Coatings, NATO Advanced Study Institute Series B: Physics, Vol. 266, Plenum, New York, 1991.
5. M.J. Mirtich, D.M. Swec, J.C. Angus, Dual Ion Beam Deposition of Carbon Films with Diamond-Like Properties, NASA TM 83743.
6. American National Standard for Ophthalmics-Nonprescription Sunglasses and Fashion Eyewear Requirements, September 25, 1986.
7. Boresi, Arthur P. and Sidebottom, Omar M., Advanced Mechanics of Materials, Fourth Edition, John Wiley & Sons, New York, 1985.
8. Bull, S.J., "Failure Modes in Scratch Adhesion Testing," Surface and Coatings Technology, 50 (1991) 25-32.
9. Mirtich, M.J., Journal Vacuum Science and Technology, 18, 2, March, 1981.

TABLE I
Potential Diamondlike Carbon Film Applications

DLC Film Application	Desirable DLC Film Properties
1. Protective coating for sunglass lenses	Hardness, scratch resistance
2. Protective coating for eyeglass lenses	Hardness, scratch resistance, transmittance
3. Hermetic coating for eyewear	Hermeticity, hardness, scratch resistance, transmittance
4. Abrasion, moisture resistant coating for optical surfaces (visible and infrared)	Hermeticity, hardness, scratch resistance, transmittance
5. Magnetic recording head	Hermeticity, hardness, scratch resistance
6. Coating computer hard disk	Hermeticity, hardness, scratch resistance
7. Abrasion resistant coating for optical windows in bar code scanners	Hardness, scratch resistance
8. Biomedical applications	Biocompatibility, hermeticity
9. Chemically resistant protective coating	Hermeticity, hardness
10. Enhanced IR transmittance of Ge and Si Infrared optics	IR transmittance, index of refraction, abrasion protection
11. Cutting blades	Smoothness, hardness, hermeticity
12. Abrasion resistant non-stick coating for cookware	Hardness, scratch resistance

TABLE II
Focused Diamondlike Film Applications
at NASA Lewis Research Center

- (1) Abrasion-Resistant Anti-Reflective Optical Coatings
 - Plastic sunglass lenses
 - Plastic eyeglass lenses
 - Other optical substrates such as glass

- (2) Chemical/environmental protection ("hermetic sealing") of transparent substrates
 - Quartz
 - Glass
 - Plastics

- (3) Wear Protection of Non-Optical Substrates
 - Magnetic Disks
 - Cutting Surfaces
 - Other Wear Parts

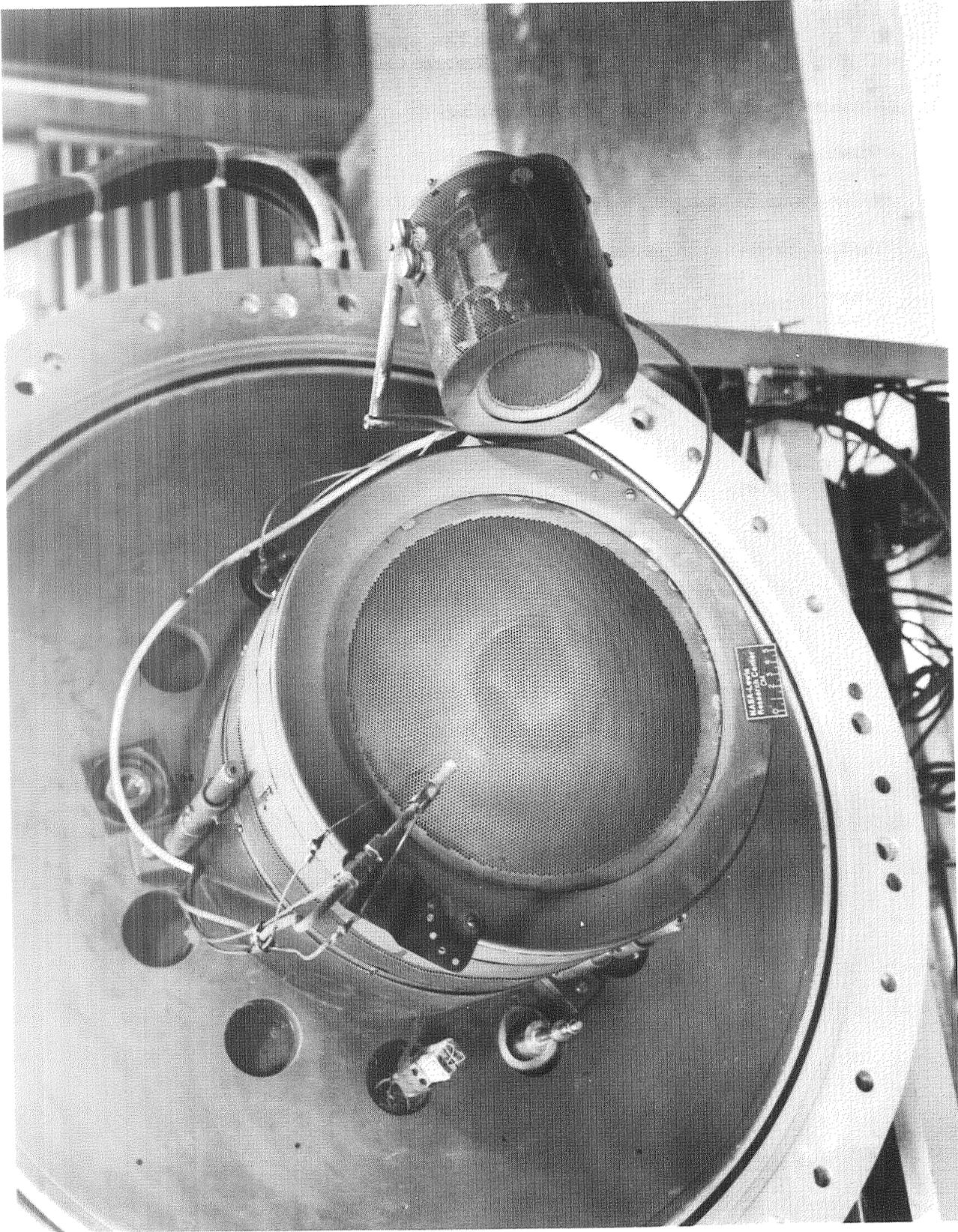


Figure 1. Dual Ion Beam Apparatus

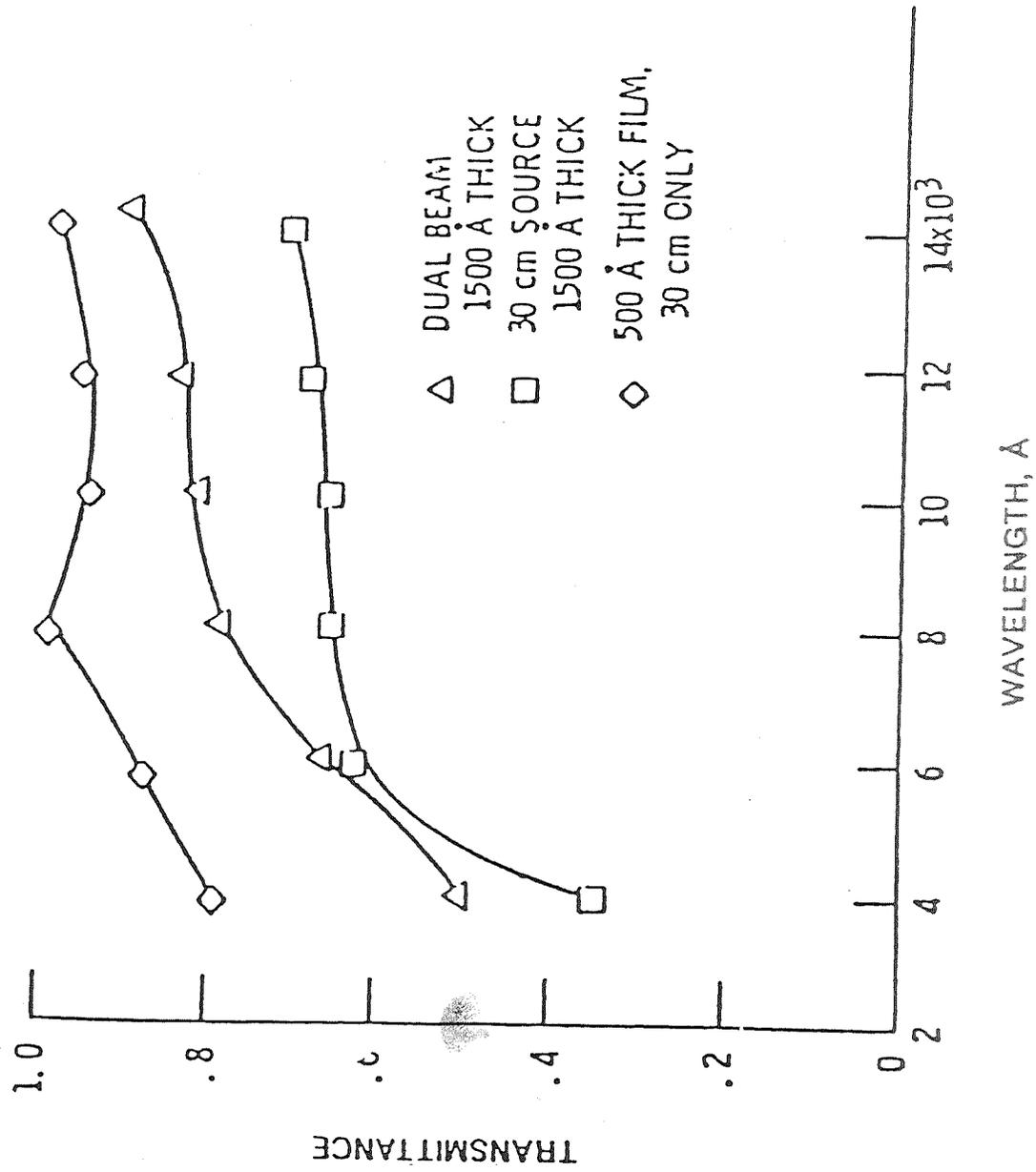


Figure 2. Transmittance versus wavelength for DLC films using CH₄ in dual beam or single ion sources.

TRANSMITTANCE AT 0.5 μm for 1000 Å Thick DLC Film

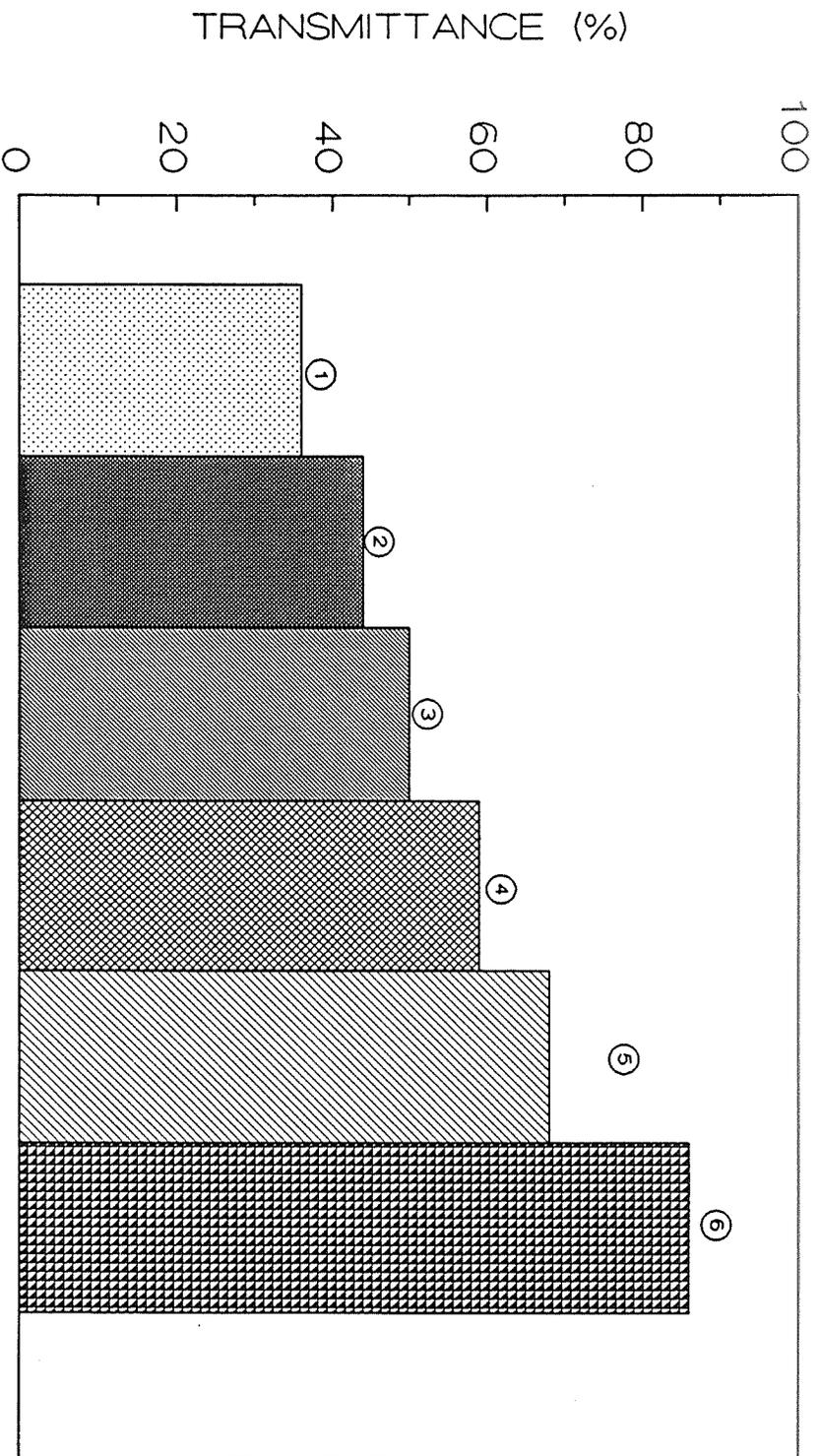
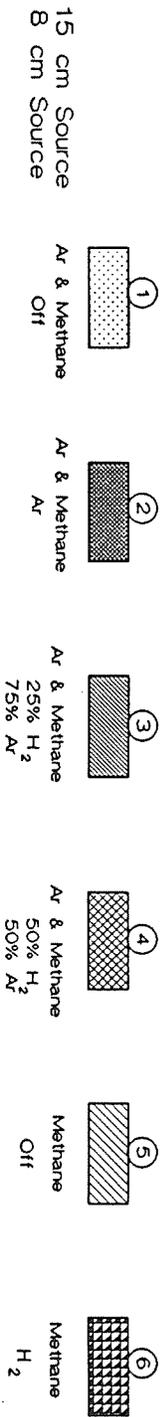


Figure 3. Transmittance at 0.5 μm for 1000 Å thick DLC films deposited on quartz under various dual beam conditions

DLC DEPOSITION CONDITIONS



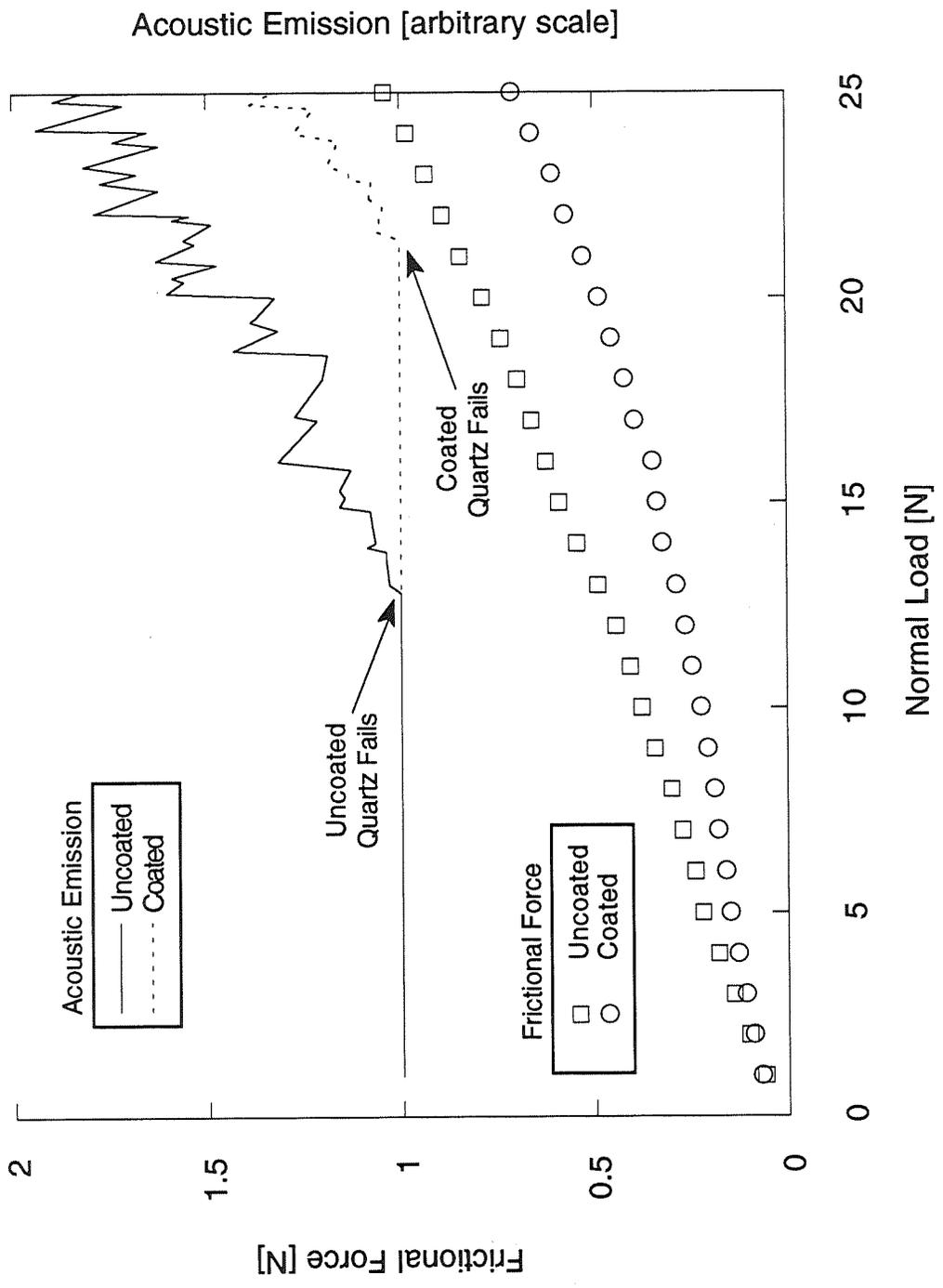
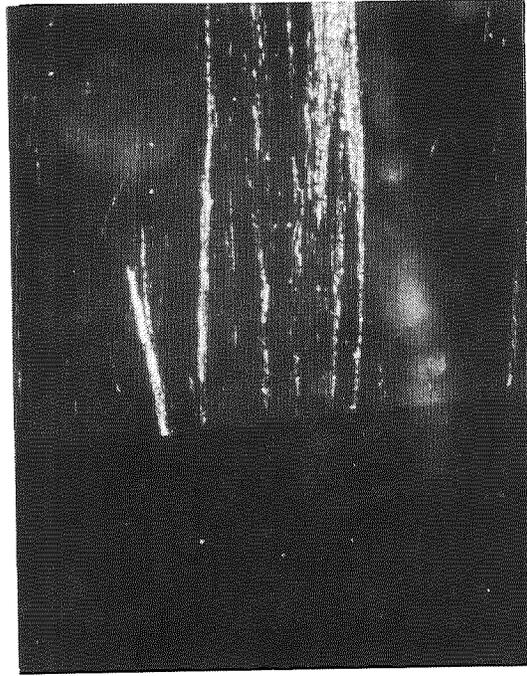


Figure 4. Coated quartz shows lower friction and withstands higher normal load before fracture (as indicated by acoustic emission).

ELECTRO-PHYSICS BRANCH



↑
DLC FILM

↑
QUARTZ

SCRATCH TEST OF DLC FILM USING SiO₂ PARTICLES

Fig.5 Scratch test of DLC film using SiO₂ particles.

CD-89-44216

PROPERTIES OF EXTRUDED PS-212 TYPE SELF-LUBRICATING MATERIALS

W.J. Waters
Materials Engineer
Sverdrup Technologies, Inc.
Cleveland, OH 44135

H.E. Sliney
Materials Engineer
NASA Lewis Research Center
Cleveland, OH 44135

R.F. Soltis
Electrical Engineer
Cortez III
Cleveland, OH 44135

N 93-25365

150474

p. 15

ABSTRACT

Research has been underway at the NASA Lewis Research Center since the 1960's to develop high temperature, self-lubricating materials. The bulk of the research has been done "in-house" by a team of researchers from the Materials Division.

A series of self-lubricating solid material systems has been developed over the years. One of the most promising is the composite material system referred to as PS-212 or PM-212. This material is a powder metallurgy product composed of metal bonded chromium carbide and two solid lubricating materials known to be self-lubricating over a wide temperature range.

NASA feels this material has a wide potential in industrial applications. Simplified processing of this material would enhance its commercial potential. Processing changes have the potential to reduce processing costs, but tribological and physical properties must not be adversely affected.

Extrusion processing has been employed in this investigation as a consolidation process for PM-212/PS-212. It has been successful in that high density bars of EX-212 (extruded PM-212) can readily be fabricated. Friction and strength data indicate these properties have been maintained or improved over the P.M. version. A range of extrusion temperatures have been investigated and tensile, friction, wear and microstructural data have been obtained. Results indicate extrusion temperatures are not critical from a densification standpoint, but other properties are temperature dependent.

INTRODUCTION

As advances in engine technology have proceeded over the years, a corresponding need has been generated to develop lubricating systems capable of performing in extreme temperature environments. These temperatures exceed the useful range of normal lubricating systems. Liquid and dry lubricants cannot operate at temperatures where the organic components break down or the solids oxidize. New methods include gas bearings and self-lubricating/load bearing composites such materials are needed to fulfill the promise of improved high temperature system performance.

A series of self-lubricating composite materials, unique in their chemistries, has been developed at NASA Lewis. The latest and most promising in this series is referred to as PS-212 or PM-212. This material is composed of powders of chromium carbide with a nickel base alloy binder combined with silver and eutectic fluorides. Both PS-212 and PM-212 have the same composition, only fabrication methods differ. The first version was PS-212, a plasma sprayed material that is described in detail in refs. 1, 2, 3, and 4. Plasma spraying of the material, while relatively simple, has some inherent drawbacks including the uneven buildup of the sprayed surface, compositional variation, entrapped porosity and oxidation, and the inability to coat many internal surfaces. PM-212 was the next version and it was developed using standard powder metallurgy techniques including cold or hot pressing, cold isostatic pressing and sintering or hot

isostatic pressing. This data is reported in refs. 5 and 6. Reference 5 deals with powder metal processing and resultant tribological properties. In ref. 6 strength properties of PM-212 are reported for various processing conditions. Thermal conductivity and thermal expansion are also included in this reference.

A new material must be cost effective to be accepted as a viable commercial product. Raw material costs and processing are major considerations in this acceptance process. Extrusion has the potential for reduced processing costs and is a commonly accepted processing procedure.

Optimum extruded properties should only be expected if the variables are optimized. These include: (1) starting powder sizes and distribution, (2) extrusion temperatures, (3) reduction ratio and extrusion speed. This investigation was limited and only looked at varying extrusion temperatures over a range from 1400°F to 1800°F. The results reported are therefore not necessarily optimum for EX-212.

The advantages for using extrusion include the following: Commercial extrusion is done over a range of temperatures and is used on a variety of materials including both metallics and non-metallics. The extrusion process both consolidates and shapes. Shapes and sizes can vary greatly depending upon tooling. Resultant products are consistent in properties and reduced in cost.

MATERIALS AND PROCEDURE

PS-212, PM-212 and EX-212 are three process versions of the same basic material. They are composed of 70% of a metal bonded, chromium carbide powder (i.e. Metco 430 NS), 15% silver, and 15% barium fluoride/calcium fluoride. The chemical compositions of these components are listed in Table I. Metco 430 NS powder gives the hardness, wear resistance and load bearing capacity while silver lubricates from cryogenic temperatures to 900°F and the fluorides lubricate from 900°F to 1600°F.

MATERIAL PREPARATION

All components were obtained from a commercial source that processed and blended the powders under the direction of NASA personnel. As indicated in Table I, Metco 430 NS powder is from -200 to +400 mesh, silver powder is -100 to +325 mesh and the eutectic fluorides are -200 to +325 mesh. These powders are mixed in a "V-blender" for 45 minutes.

The blended powders are processed in air and are gravity filled into the extrusion cans. These cans are fabricated from mild steel and have a finish size of 3" O.D. x 2" I.D. x 6" long. After filling the cans, the contents are degassed prior to welding on the top. After sealing the cans were furnace heated for one hour at the extrusion temperature and then rapidly transferred to the extrusion press. Area reductions for the extrusions were held constant at 16:1. This resulted in an EX-212 bar diameter of about .4 inch and a length of 5 feet. The extrusions were done in a 1020 ton vertical press. Maximum punch pressure is less than 190,000 PSI. Extrusion temperatures were 1400°F, 1500°F, 1600°F, 1700°F, and 1800°F.

After extrusion, the bars were inspected by both x-ray and surface fluorescence (Zygló) and then fabricated into test samples. Two types of test samples were used in this investigation. They are shown in fig. 1 and include tensile test bars (ref. 7) and friction and wear bars. Test specimens were fabricated by diamond grinding for the friction wear bars and by silicon carbide grinding for the tensile bars. Both methods worked, but the diamond ground surface was smoother.

TENSILE AND TENSILE TESTING

Tensile tests were run at room temperature, 800°F, 1200°F, 1400°F and 1600°F. They were run, in most cases, as single data points. Tests were run in a controlled strain rate test machine and temperatures were measured with thermocouples attached to the test section.

Friction and wear specimens were scrubbed with an alumina/water paste and rinsed with deionized water prior to testing. These specimens were slid against a Rene '41 (nickel base alloy) counterface disk. Testing was done in a pin-on-disk tribometer. The hemispherically tipped (3/16" radius of curvature) specimen (pin) was loaded against the flat surface of the rotating Rene '41 disk. Sliding speed was 9 ft/sec with a deadweight load of 1.1 lb. Tests were run in an air atmosphere (relative humidity 35% at room temperature) at 70°F, 660°F and 1400°F.

Nine double-ended EX-212 specimens were tested for a total of 36 half hour tests by running two tests on each pin end. Of the nine specimens, three each were from the 1400°F, 1600°F and 1800°F extrusions. The Rene '41 disks were heated by induction and were generally run at 1000 rpm (9.2 ft/sec.). (A velocity survey was conducted by running successive five minute tests at 10, 100, 500, 1000, 2000, and 3500 rpm.)

After each half-hour test, wear measurements were made. The pin tips were photomicrographed to determine wear scar diameters. A surface profile meter was used to measure the cross sectional area of the disk wear track. These measurements were used to calculate wear volumes and wear factors for the pins and disks.

Microstructures of the extruded EX-212 material were examined in both the etched and the unetched conditions. All microstructural examination was done with optical microscopes on longitudinal sections. It should be noted that optical viewing shows the fluorides as black, due to their low reflectivity. They should not be mistaken for voids which have no reflectivity.

HARDNESS

Hardness tests were run on an automatic readout Rockwell Hardness Tester. All tests were run on cross sectional segments cut from the extrusions. Tests were run on the R "C" scale and are listed in Table IV.

RESULTS AND DISCUSSION

Extrusions of PS-212 blended powder were run at temperatures of 1400°F, 1500°F, 1600°F, 1700°F and 1800°F. This was done to determine the feasibility of using a low cost commercial process to consolidate the composite material into near net shape. Starting powder sizes were those used in the earlier studies of PS/PM-212. Reduction ratios were similar to those used in earlier studies with high temperature alloy systems. (Ref. 7)

All of the extrusions were successful in that they did consolidate into a bar configuration. X-rays of the bars indicated a relatively uniform, continuous extrusion with no bursts or discontinuities. Test samples were examined by X-ray and Zyglo and found to be generally porous. The porosity was the greatest in the lower temperature extrusions. Less porosity was evident as extrusion temperature increased and was a minimum at 1800°F.

TENSILE

Tensile data is shown in Table II and plotted in fig. 2. Generally the lower temperature extrusions had better strength at test temperatures from room to 1200°F and the higher temperature extrusions had better strengths at 1400°F and 1600°F. Elongation in all cases were less than 3%. The low temperature extrusions had the least ductility as measured by tensile elongation. The brittle fluoride eutetic encapsulates the metallic phases and limit both the strength and ductility. Both the strengths and the ductility of the EX-212 compare favorably with the PM-212. It should also be noted that refinements in the particular size of the components may have a strengthening effect on the resultant extrusion, especially at the lower extrusion temperatures.

FRICITION AND WEAR

The friction and wear results are shown in Table III and figs. 3-8. The friction and wear behavior for the EX-212 samples are shown to be similar to results previously found for sintered PM-212.

The friction values for the EX-212 at a sliding velocity of 9.2 ft/sec ranges from a low of 0.26 at 1400°F for the 1400°F extrusion to a high of 0.50 at 660°F for the 1600°F extrusion. The pin wear factors for all the extruded tests were determined to fall within the moderate to low range within the 1400°F extrusion showing the least wear at 1400°F. The disk wear factors were found to fall within the moderate to low range for all the extruded tests except two. The wear factor for the 1600°F and 1800°F extrusions were negative, denoting material transfer from pin to disk. This transfer was reflected in the wear factors of the

corresponding pins, which had the two highest pin wear volumes of the extruded tests.

In general for the EX-212, increasing the disk speed was found to decrease the friction coefficient, with less of an effect at higher temperatures. This behavior was also displayed by the sintered PM-212 previously studied.

Overall the EX-212 samples were determined to display friction and wear behavior to sintered PM-212. Varying the extrusion temperature was found to have little effect on the tribobehavior of EX-212, although the 1400°F extrusion has the best overall performance.

MICROSTRUCTURE

Fig. 9 shows 25X unetched photomicrographs of all extrusion temperature bars. Stringers appear most evident in the lower temperature bars. Some voids are present in all bars but again they are the most evident in the low temperature extrusions. Fig. 10 shows the etched bars at 100X. Particles are deformed and elongated (cold worked) in the 1400°F extruded bars. This deformation decreases as the extrusion temperature is increased. At 1800°F the extruded microstructure is equiaxed. Fig. 11 shows the etched microstructure at 1000X. The 1400°F microstructure shows three major components. Fluorides are shown as dark phases, the white phase is silver and the grey phase is the Metco 430NS. Note that there are very thin bands of silver between grey phases. As extrusion temperature is increased, a second phase appears in the matrix of the Metco 430NS. The carbides remain clear. Silver appears to coalesce as the extrusion temperature increases.

It should be emphasized that the interconnected dark grey areas are fluorides and not voids. Density measurements of the extrusions indicate near theoretical density in all extrusions.

HARDNESS

The hardness tests were taken on the Rockwell "C" scale and hence represent an average hardness over a range of phases. Micro hardness studies might be appropriate for further studies on the effect of extrusion parameters. The hardness values obtained are listed in Table IV and generally indicate an increasing hardness with increasing extrusion temperatures. This is unexpected since lower temperature working conditions would be expected to result in increased work hardening and hardness. One explanation might be the increased evidence of carbide growth in the higher temperature extrusions and the development of the second phase precipitate in the nickel matrix.

CONCLUSIONS

The PS-212 chemistry and the components lend themselves to fabrication by the extrusion process. Large reduction ratios (16:1) of the canned powders result in solid bars with minimal porosity. Component particle size is probably a factor in resulting physical properties. Properties of the extruded EX-212 compare favorably with the prior versions of the PM-212.

Extrusion offers a low cost process to fabricate near net shape components for bearing applications. Further work might be done to evaluate finer particulate sizes and their effect on resultant physical properties.

REFERENCES

1. Sliney, H.E.: Carbide/Fluoride/Silver Self-Lubricating Composite; U. S. Patent 4,728,448. Useful at Low and High Temperatures. March 1, 1988.
2. Sliney, H. E. and Della Corte, C.: Method of Making Carbide/Fluoride/Silver Composites. U. S. Patent Application 571,058, August 1990.
3. Della Corte, C. and Sliney, H. E.: Composition Optimization of Self- Lubricating Chromium-Carbide-Based Composite Coatings for Use to 760 C. ASLA Trans., Volume 30, No. 1, January 1987, pp. 77-83.
4. Sliney, H. E.: Hot Piston Ring/Cylinder Liner Materials Selection and Evaluation. SAE Paper 880544, 1989 (also NASA TM 100276).
5. Della Corte, C. and Sliney, H. E.: Tribological Properties of PM-212: A High Temperature Self-Lubricating, Powder Metallurgy Composite, NASA TM 102355.
6. Edwards, P. M. and Sliney, H. E.: Mechanical Strength and Thermophysical Properties of PM-212: A High Temperature Self- Lubricating Powder Metallurgy Composite, NASA TM 103694.
7. Waters, W. J. and Freche, J. C.: Strength Enhancement Process for Prealloyed Powder Superalloys, NASA TM 78834.

TABLE I. - COMPONENTS OF EX-212

Component	Composition, wt %	Composition, vol %	Particle size U.S. Sieve No. (μm)	MP ¹ °C	Hardness, Hv kg/mm ²
Component A: Bonded Chromium Carbide: 70 wt % of PH212					
Cr ₂ C ₃	45	47	-200 + 400 (35 to 74)	1895	^a 1300
Ni	28	22		1455	^a 570
Co	12	10		-----	-----
Cr	9	9		-----	-----
Mo	2	1		-----	-----
Al	2	5		-----	-----
B	1	3		-----	-----
Si	1	3		-----	-----
Component B: Silver Metal: 15 wt % of PH212					
Ag	100	100	-100 + 325 (44 to 150)	961	^c 25
Component C: Prefused Eutectic: 15 wt % of PH212					
BaF ₂	62	52	-200 + 325 (44 to 74)	1050 1280	^b 150 ^b 110
CaF ₂	38	48		1423	^b 145

^(a) Handbook of Physics and Chemistry, 70th ed., 1990, CRC Press Inc., Boca Raton, FL.

^(b) Deadmore, D.L. and Sliney, H.E.: "Hardness of CaF₂ and BaF₂ Solid Lubricants at 25 to 670 °C," NASA TM-88979, March 1987.

^(c) Lozinskii, M.G.: "High Temperature Metallography," Pergamon Press, 1961.

TABLE II TENSILE STRENGTHS OF EX-212

Extrusion Temp, Deg. F	Test Temp., Deg. F	Ultimate Tensile Strength, KSI
1400	70	21.4, 23.4
	800	20.0
	1200	17.5
	1400	7.0
	1600	1.6
1500	70	25.5
	800	22.0
	1200	17.0
	1400	7.1
	1600	2.2
1600	70	24.4, 24.0
	800	20.0
	1200	13.0
	1400	7.0
	1600	1.6
1700	70	19.0, 19.5
	800	16.5
	1200	12.0
	1400	5.6
	1600	1.3, 1.2
1800	70	24.2, 26.5
	800	19.0
	1200	16.0
	1400	10.5
	1600	2.4

TABLE III - Friction and Wear Comparison of EX-212 Pins Against René 41 Disks With Sliding Velocity of 9.2 ft/sec at 1000 RPM

Processing Method	Temperature [°C]	Friction Coefficient	K _{pin} 10 ⁻⁵ mm ³ /N-m	K _{disk} 10 ⁻⁵ mm ³ /N-m	Number of Tests
Extruded at 1400F	25 (77F)	0.33±0.01	2.1±0.5	0.55±0.4	3
	350 (662F)	0.46±0.07	1.5±1.3	0.23±0.8	3
	760 (1400F)	0.26±0.03	0.39±0.1	0.15±0.5	3
Extruded at 1600F	25	0.37±0.04	3.2±0.3	3.5±4.3	3
	350	0.50±0.05	1.6±0.6	0.18±0.2	3
	760	0.35±0.01	6.8±3.7	-3.0±2.1	3
Extruded at 1800F	25	0.35±0.04	3.1±2.3	1.4±3.2	3
	350	0.44±0.08	1.2±0.6	0.23±1.1	3
	760	0.33±0.03	3.5±9.4	-1.1±3.9	3
Sintered	25	0.35±0.05	3.2±1.5	7.0±2.0	2
	350	0.38±0.02	3.9±1.8	0.35±0.1	3
	760	0.35±0.06	0.36±0.09	1.0±6.0	4
HIPped	25	0.37±0.04	1.8±0.4	0.45±0.1	≥4
	350	0.32±0.07	2.5±0.3	0.85±0.4	≥4
	760	0.31±0.04	0.07	2.2±0.8	≥4

Note: Uncertainties represent one standard deviation from the mean for the friction coefficients and the data scatter range for the wear factors.

**TABLE IV
ROOM TEMPERATURE HARDNESS DATA FOR EX-212**

EXTRUSION TEMP, DEG F	AVERAGE HARDNESS, R"C"
1400	21.5
1500	29.0
1600	31.5
1700	24.5
1800	32.0

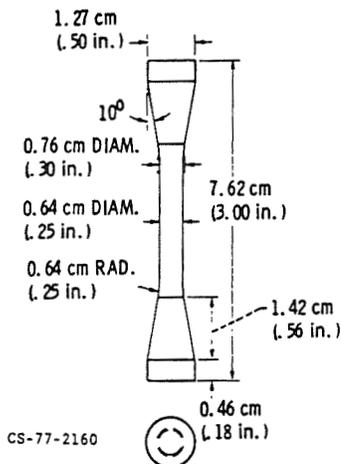


FIGURE 1.A - Tensile/Stress Rupture Specimen Ground From Extruded Bar Stock



****Finishing of Hemispherical Tips****

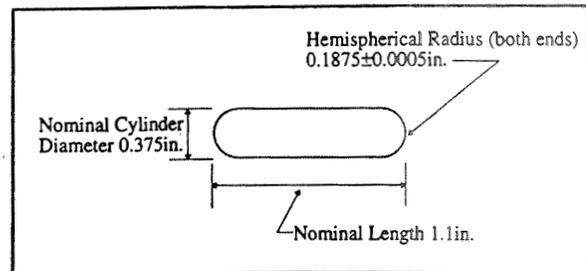
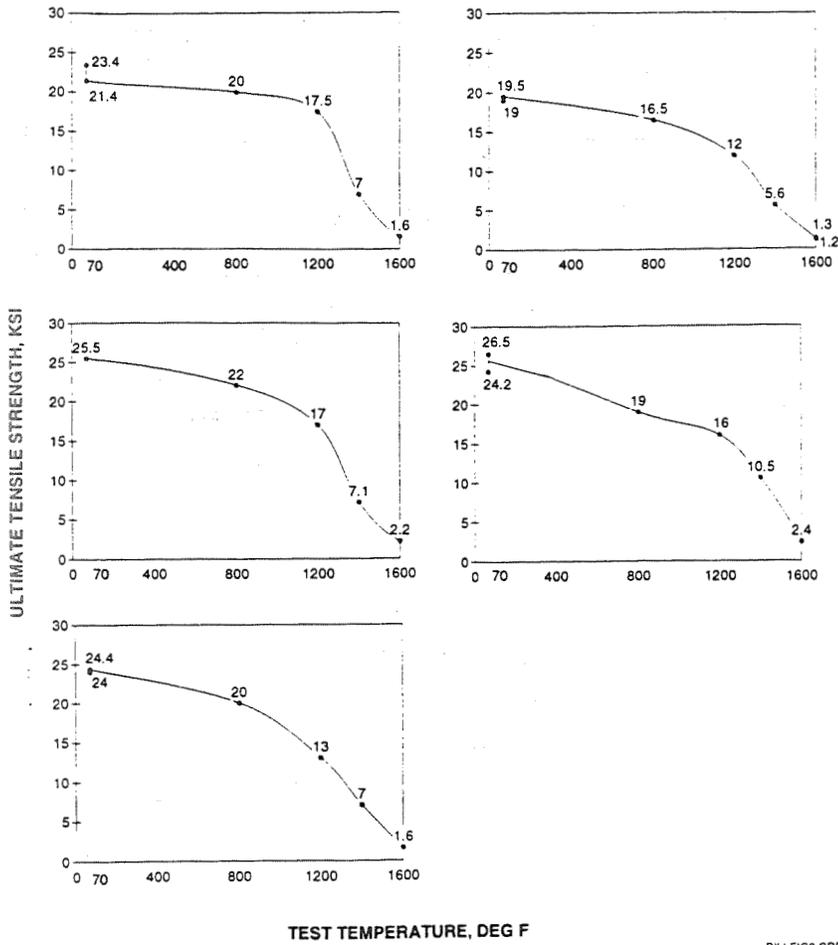


Figure 1.B

FIG. 2 ULTIMATE TENSILE STRENGTH OF EX-212 AT VARIOUS TEMPERATURES



BILLFIG2.GRF

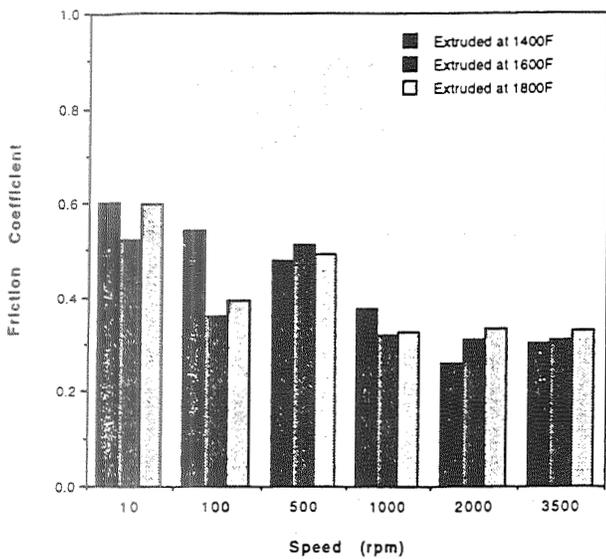


Figure 3 - Velocity Survey for Extruded PM212
Test Temperature at 25C (77F)

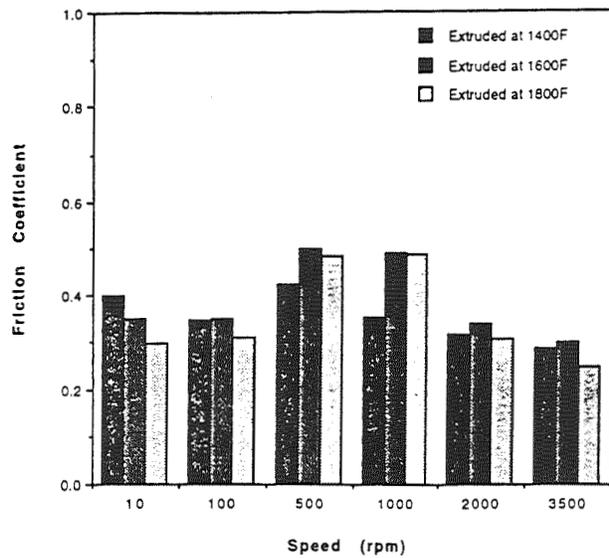


Figure 4 - Velocity Survey for Extruded PM212
Test Temperature at 350C (662F)

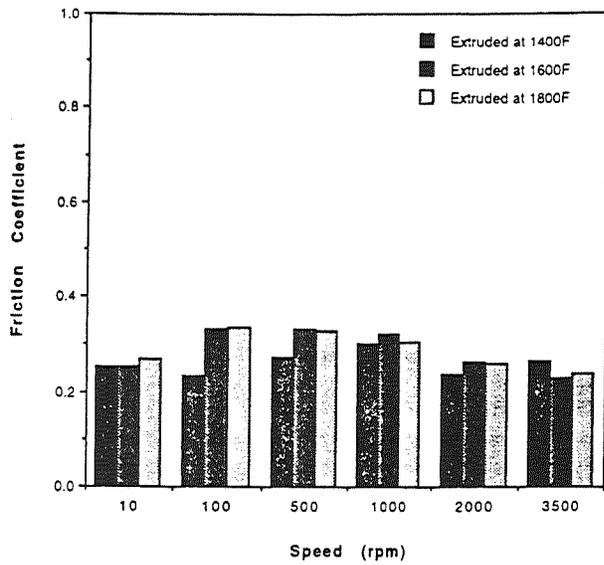


Figure 5 - Velocity Survey for Extruded PM212
Test Temperature at 760C (1400F)

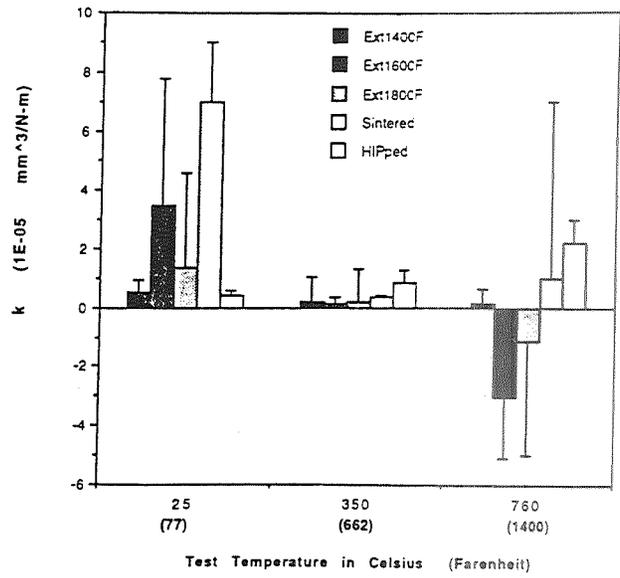


Figure 6 - Disk Wear Factors Rene 41 Disks vs.
PM212 Pins

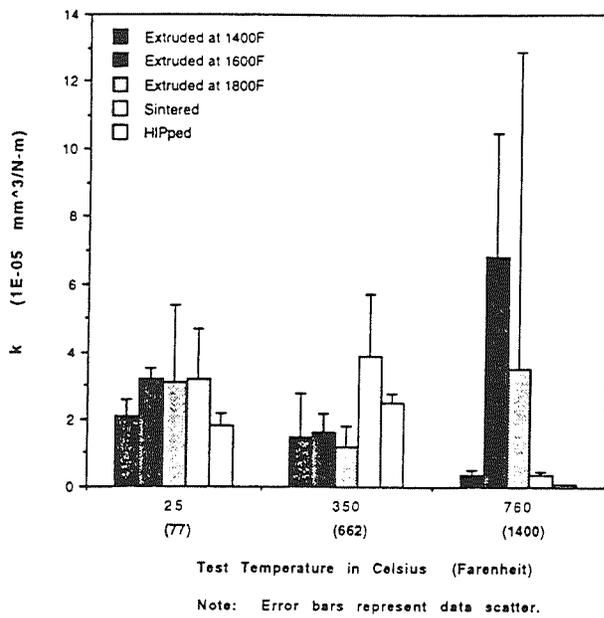


Figure 7 - Pin Wear Factors PM212
Comparison

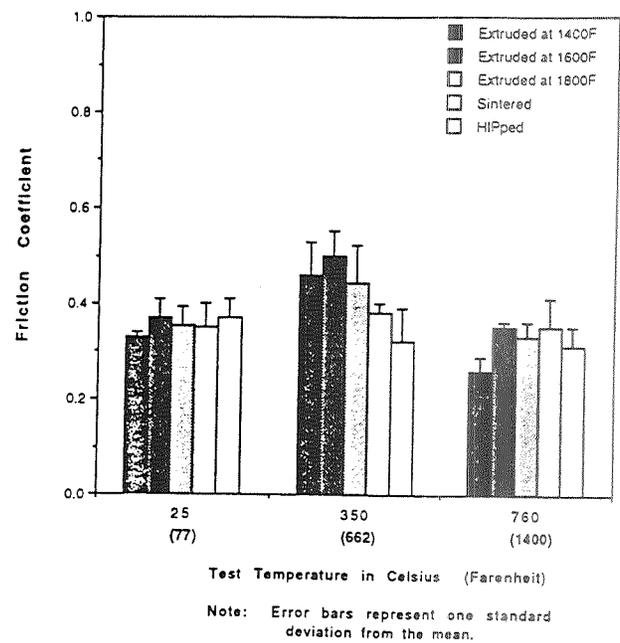
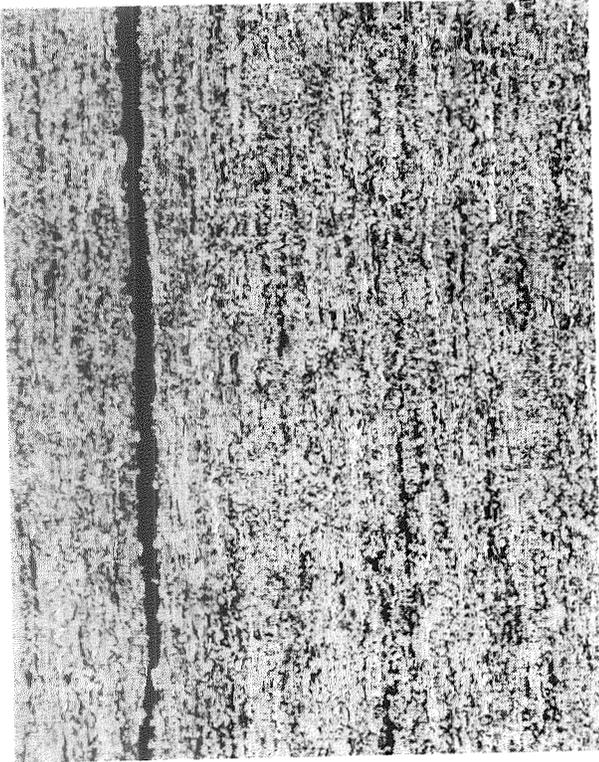


Figure 8 - Friction Comparison for PM212

(a) Ext. at 1400°F



(b) Ext. at 1500°F

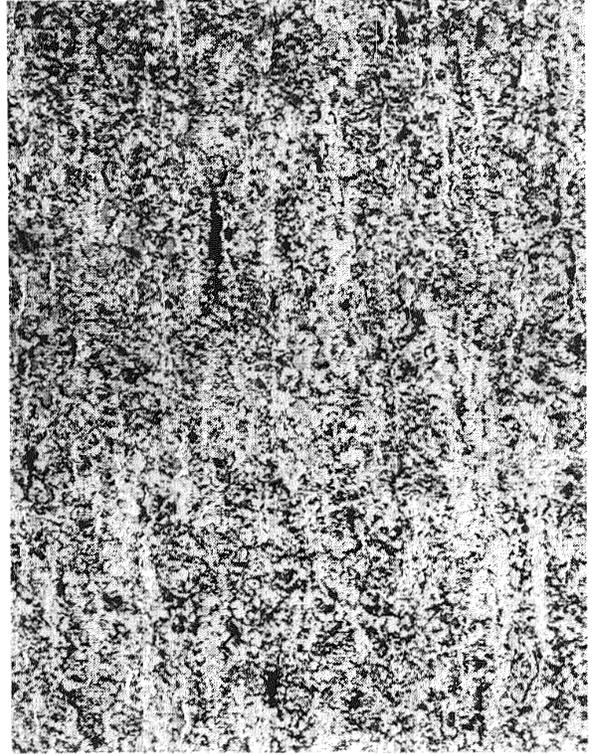


FIG. 9 EX-212 Unetched, 25x

(c) Ext. at 1600°F



(d) Ext. at 1700°F

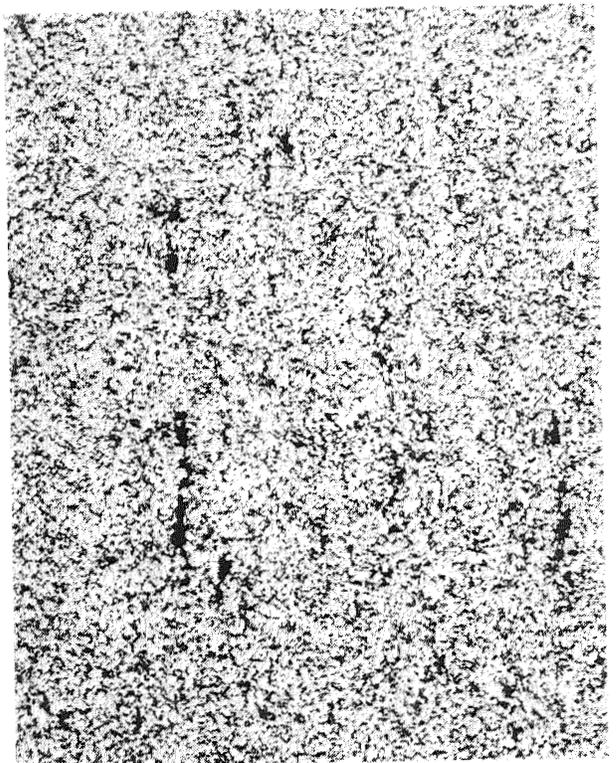
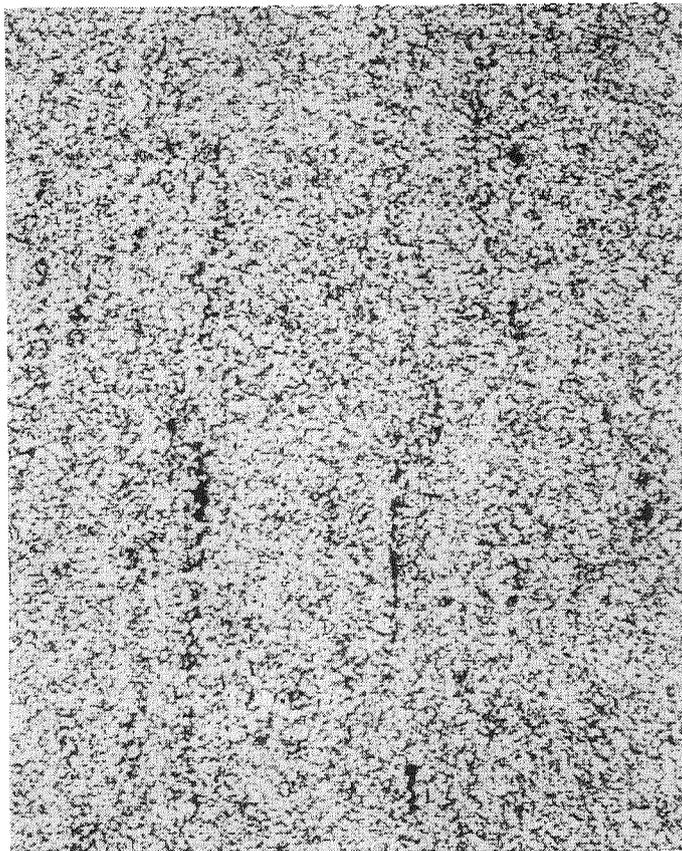
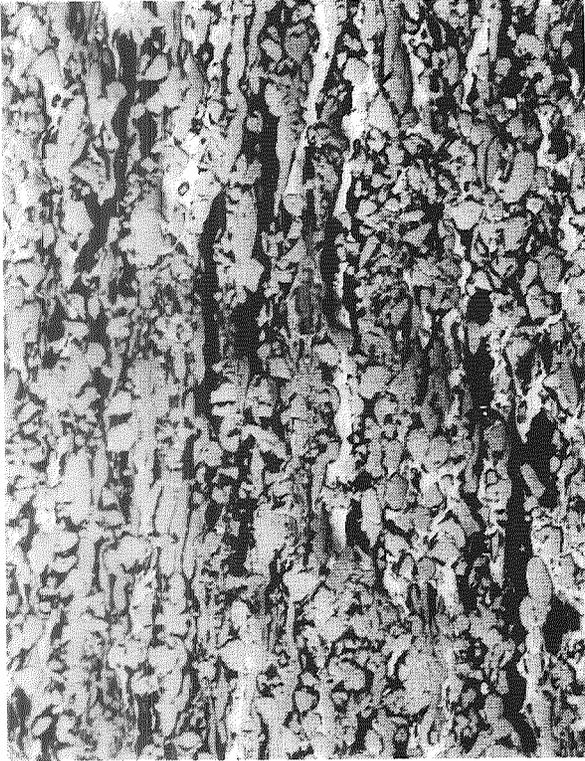


FIG. 9
(e) Ext. at 1800°F



(a) Ext. at 1400°F

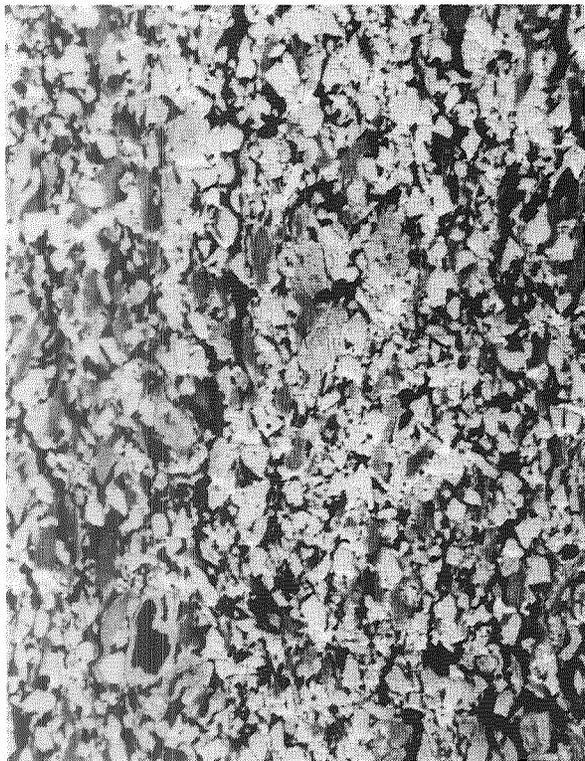


(b) Ext. at 1500°F



FIG 10. EX-212 Etched, 100x

(c) Ext. at 1600°F



(d) Ext. at 1700°F

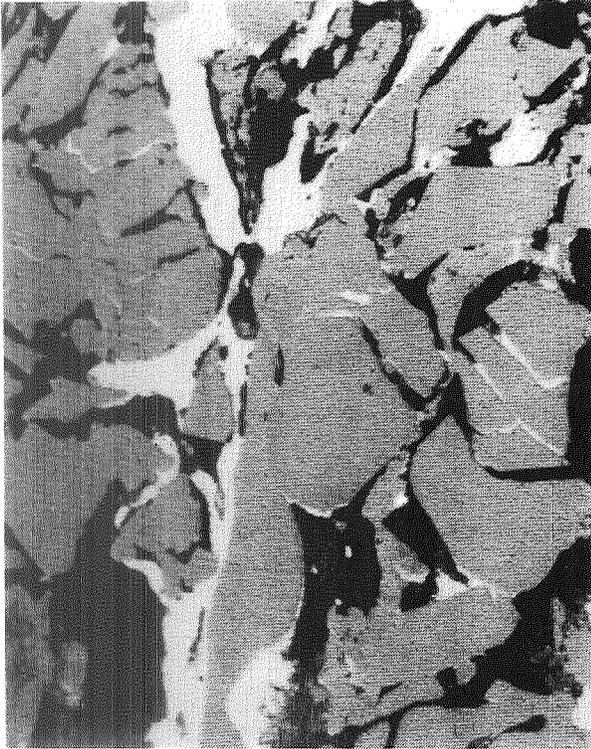


FIG. 10

(e) Ext. at 1800°F



(a) Ext. at 1400°F



(b) Ext. at 1500°F

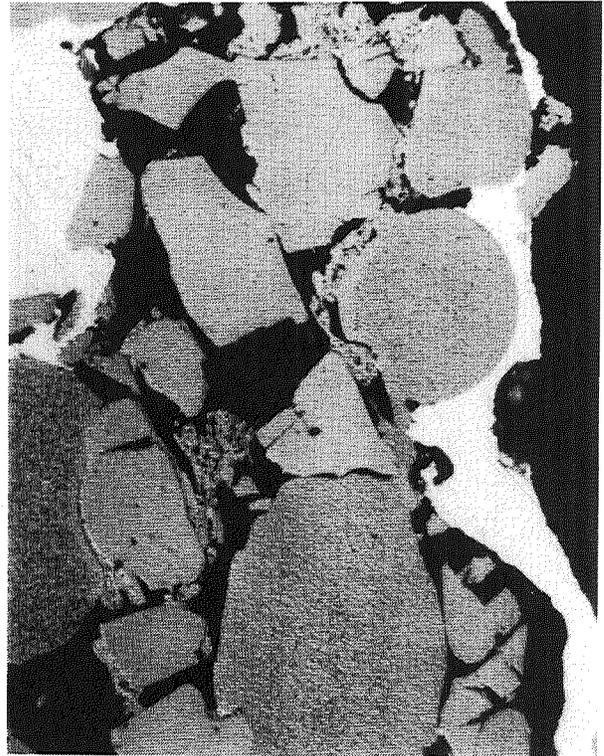


FIG. 11 EX-212 Etched, 1000x

(c) Ext. at 1600°F



(d) Ext. at 1700°F

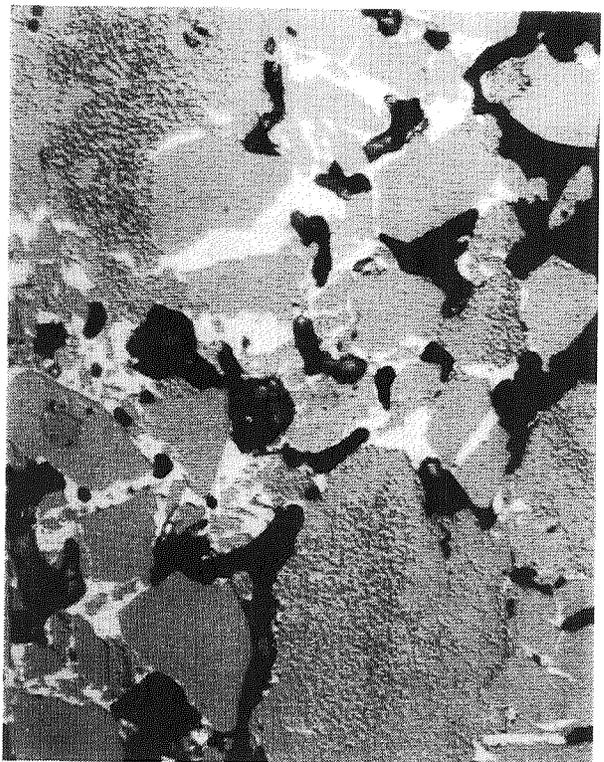
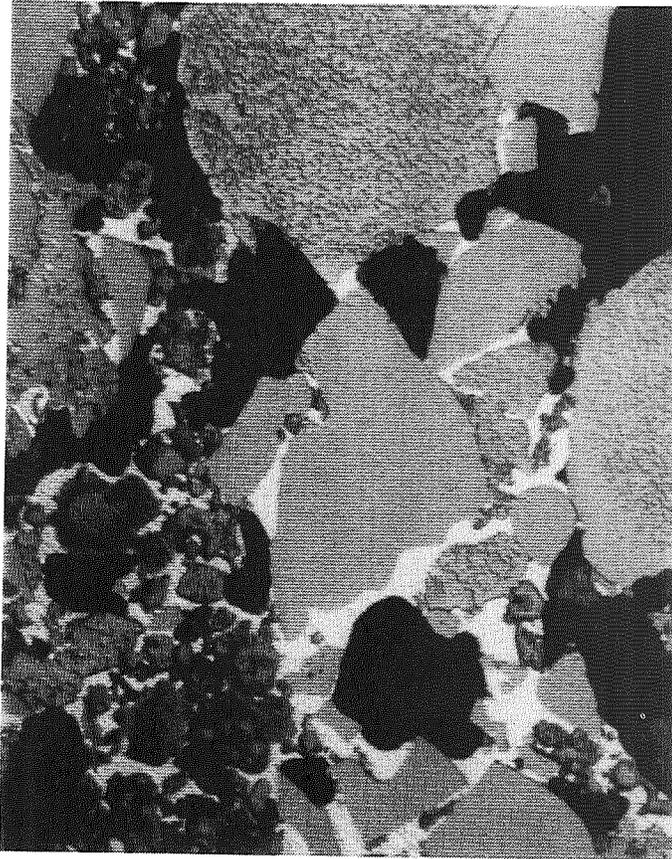


FIG. 11

(e) Ext. at 1800°F





**BIOTECHNOLOGY AND LIFE SCIENCES
PART 1**

PRECEDING PAGE BLANK NOT FILMED

57.

56
~~INTENTIONALLY BLANK~~



55-51
150475
p. 12

N93-25560

"MEASURING THE METASTATIC POTENTIAL OF CANCER CELLS"

Dennis R. Morrison, Ph.D.
NASA - Johnson Space Center
Houston, Texas 77058

Howard Gratzner, Ph.D.
DNA Sciences, Inc.
The Woodlands, Texas 77380
and

M. Z. Atassi, Ph.D.
Baylor College of Medicine
Houston, Texas 77030.

ABSTRACT

Cancer cells must secrete proteolytic enzymes to invade adjacent tissues and migrate to a new metastatic site. Urokinase (uPA) is a key enzyme related to metastasis in cancers of the lung, colon, gastric, uterine, breast, brain and malignant melanoma. A NASA technology utilization project has combined fluorescence microscopy, image analysis and flow cytometry, using fluorescent dyes, and urokinase-specific antibodies to measure uPA and abnormal DNA levels (related to cancer cell proliferation) inside the cancer cells. The project is focused on developing quantitative measurements to determine if a patient's tumor cells are actively metastasizing. If a significant number of tumor cells contain large amounts of uPA (esp. membrane-bound) then the post-surgical chemotherapy or radiotherapy can be targeted for metastatic cells that have already left the primary tumor. These analytical methods have been applied to a retrospective study of biopsy tissues from 150 node negative, stage I breast cancer patients. Cytopathology and image analysis has shown that uPA is present in high levels in many breast cancer cells, but not found in normal breast. Significant amounts of uPA also have been measured in glioma cell lines cultured from brain tumors. Commercial applications include new diagnostic tests for metastatic cells, in different cancers, which are being developed with a company that provides a medical testing service using flow cytometry for DNA analysis and hormone receptors on tumor cells from patient biopsies. This research also may provide the basis for developing a new "magic bullet" treatment against metastasis using chemotherapeutic drugs or radioisotopes attached to urokinase-specific monoclonal antibodies that will only bind to metastatic cells.

INTRODUCTION

Malignant cells are characterized by abnormal levels of DNA, rapid proliferation, uncontrolled growth and the ability to invade surrounding normal tissues. The measurement of biochemical markers on cancer cells can provide valuable information as to disease-free survival, time to relapse and thus provides the physician valuable with data for planning adjuvant therapy. Indirect immunoassays of markers extracted from biopsy tissues are important, but more precise measurements can be made by analytical cytometry. The current trend is towards microscopic analysis of the immunochemically stained tumor sections or dissociated cells, coupled with quantitation by image analysis. Specific markers can be directly associated with the cancer tissue, as opposed to biochemical extraction procedures. Tumor markers currently assessed include those which measure cellular proliferation, the presence of specific oncogenes, tumor-suppressor molecules, and cancer related proteins (see Table 1). Tumor related proteins include proteolytic enzymes which are correlated with recurrent disease and metastasis. These enzymes are involved in a cascade of proteolytic interactions with other enzymes and inhibitors which often culminate in the dispersal of invasive cancer cells through surrounding basement membranes and vascular systems and thereby allow them to relocate at metastatic sites distant from the primary tumor. Among these proteases are the plasminogen activators, their receptors and inhibitors, which together mediate key steps in the metastatic process.

58

Table 1. Examples of breast cancer prognostic markers currently used for patient assessment.

PROGNOSTIC	MARKER	DESCRIPTION
DNA CONTENT	Propidium iodide	DNA content variation from normal diploid (Aneuploidy)
PROLIFERATION	% S PHASE	% of cells undergoing DNA replication
	BrdU, IdU	DNA synthesis rate in S phase cells
	Ki 67	Cycling (dividing) cells
	PCNA	Proliferating cell nuclear antigen expressed in G ₁ , S and G ₂ phases of cell cycle
RECEPTORS	Estrogen (ER)	ER negative tumors do not respond to ER hormone therapy
	Progesterone (PgR)	PgR negative tumors indicate better disease-free survival (Stage II and beyond)
	HER 2/neu, c-myc	Oncogenes amplified or overexpressed in breast cancer
	EGFR	Epidermal growth factor Receptor, overexpressed in breast cancer
ENZYMES	Urokinase (uPA)	Plasminogen activator --> plasmin -->activates proteases
	Collagenase IV	Metalloproteinase that dissolves collagen & laminin
	Cathepsin B & D	Estrogen-related lysosomal enzymes

Abnormal DNA

The quantity of DNA in normal cells is a precise amount depending on the phases of the cell cycle. DNA can be measured by labeling with DNA specific fluorescent dyes (propidium iodide). The amount of dye (fluorescence) measured is directly proportional to the amount of DNA present. Also the cells can be exposed to DNA precursors (BrdU, IdU), then fluorescent labeled antibodies, specific for the DNA precursor, can be used to localize which cells are synthesizing DNA and how much [1]. Antibodies against proliferating cell nuclear antigen (PCNA) can also be used (see below). Fluorescent labeled cells then can be analysed in a laser flow cytometer or fluorescent microscope. A histogram of DNA content in normal cells shows a single diploid peak (at G₁ phase) and a tetraploid peak (at G₂+M phase). However, in most biopsies the abnormal DNA content of tumor cells is detected as a second G₁ peak or multiple peaks. Abnormal DNA (DNA aneuploidy) is considered as an independent indicator of tumor aggressiveness and poor prognosis that is used to supplement cytopathology grading of the tumor.

Proliferation

Flow cytometric measurement of the percentage of proliferating tumor cells that are involved in synthesizing DNA (S-phase cells) also is an independent indicator of malignancy. High percentages (15 -20%) of S-phase tumor cells usually indicates an aggressive malignancy and usually correlates well with abnormally high DNA content. The labeling index (LI) obtained by pulse-labeling cells with DNA precursors represents the rate that DNA is being synthesized in tumor cells. Usually, a LI > 4% is associated with a higher probability of recurrent malignancy [2]. Antibodies against Ki 67 and PCNA have been used as a measure of tumor cell proliferation. PCNA (also called cyclin) is an auxiliary protein of DNA polymerase-alpha [3]. PCNA normally appears in only trace amounts in G₁ and increases to maximum in S-phase then declines in G₂+M phase. In tumors, high levels of PCNA are expressed in the proliferating cells in all cell phases, whereas BrdU only labels cells in S phase [4].

Receptors

Hormone receptor density on cancer cells is often important as a marker for aggressive tumors and provides strategic information for post-surgical adjuvant therapy. In the case of breast cancer, the most common prognostic indicator for the past decade has been the number of lymph nodes that the primary cancer has spread into. More recently, hormone receptor density on tumor cells has gained importance. The lack of estrogen receptor in Stage I breast cancer has become an important predictor of earlier recurrence and poor survival. In stage II, the measurement of progesterone receptors is more important than estrogen receptors for predicting disease-free survival. There is a strong correlation between tumor receptor content, % S-phase cells and DNA aneuploidy. High proliferative activity is usually inversely related to estrogen receptor levels [6].

The protooncogenes HER-2/neu (also called erbB-2) and c-myc have normal roles in the control of cell growth and differentiation, but these are amplified and overexpressed in adenocarcinomas, lung, ovarian and breast cancer. The HER-2/neu protein appears to function as a receptor for mediators of growth and differentiation. The HER-2/neu protein has structural similarity to epidermal growth factor (EGF) which is a potent cellular mitogen. The measurement of cell surface receptors for HER-2/neu and EGF also has become important as a marker for invasive cancers and poor survival. Antibodies against HER-2/neu have been shown to arrest growth of tumor cells at late S or early G₂ phase [7].

Proteolytic Enzymes and Metastasis

Cancer cells must secrete proteolytic enzymes to dissolve the basement membranes and intracellular matrix between the densely packed normal cells in order to leave the primary tumor and migrate to a new metastatic site via the blood or lymphatic circulatory systems. Serine proteases such as plasminogen activator enzymes have been linked with the invasion of tumor cells into adjacent normal tissues and with metastasis. Urokinase is not produced in most normal cells, except for low levels in certain types of normal kidney cells, colon, gastric mucosa, and endothelial cells lining small arteries. However, urokinase is produced in many tumors such as breast [8, 9], lung [10], colon [11], gastric mucosa [12], uterine [13], bladder [14], prostate [15], and malignant melanoma [16]. Both urokinase-type plasminogen activator (uPA) and tissue-type plasminogen activator (tPA) enzymes have been studied using assays of the enzymes after extraction from tumor cells or assays of supernatant medium from tissue culture of the tumor cells [12, 13, 14, 15]. In most tumors, high levels of uPA, not tPA, has been correlated with metastasis or recurrent disease [13, 14].

High levels of urokinase (>3.49 ng/mg of total protein) extracted from breast tumor tissues have recently been shown to be a good prognostic indicator for high risk of recurrence and shorter patient survival times [17]. Primary lung and colon tumor cells also produce more uPA than metastatic cells, but different methods of extraction and assays often give widely variable results [13]. Total uPA measured from tumor tissue or secreted by cultured explants is difficult to quantitate, especially if the measurements are made on a large group of cells. The data obtained is an average value of all normal and cancer cells, rather than a measurement of each individual cell. Few direct measurements of intracellular and extracellular urokinase have been made [10,18]. Urokinase (uPA) can be present in the tissues in several molecular forms. The inactive proenzyme is a single chain protein (scuPA) that is cleaved at Lys.158 to form the double chain, high molecular weight active form (HMW-uPA) that is 54 kDaltons. A low molecular weight form (LMW-uPA) can also be formed by cleavage of the HMW-uPA at Lys.135 - Lys.136 giving a 35 kD active enzyme. The active urokinase enzyme converts plasminogen into plasmin, which in turn, dissolves intracellular fibrin matrix components as well as activating collagenases, laminases, and other related protease enzymes which are important to the anchorage and growth regulation of cells (see Figure 1). Recently, it has been shown that the HMW active form of urokinase, bound to the tumor cell membrane, is responsible for the local lysis of the extracellular matrix, hence the tissue invasion mechanism for metastasis [10]. Receptor (membrane) bound uPA is twice as efficient (catalytically) as free fluid-phase uPA [19]. The unbound uPA and the LMW form is not responsible for most of the local dissolution of extracellular matrix in the immediate vicinity of the metastatic tumor cell. The presence of plasminogen activator inhibitors (PAI-1, PAI-2) also are correlated with a better prognosis and are inversely related to high levels of uPA (poor prognosis). PAI-1 binds to the active HMW-uPA, but not to the inactive scuPA [20]. Also after PAI-1 binds to the membrane-bound active uPA the complex is internalized into the cell and degraded [21].

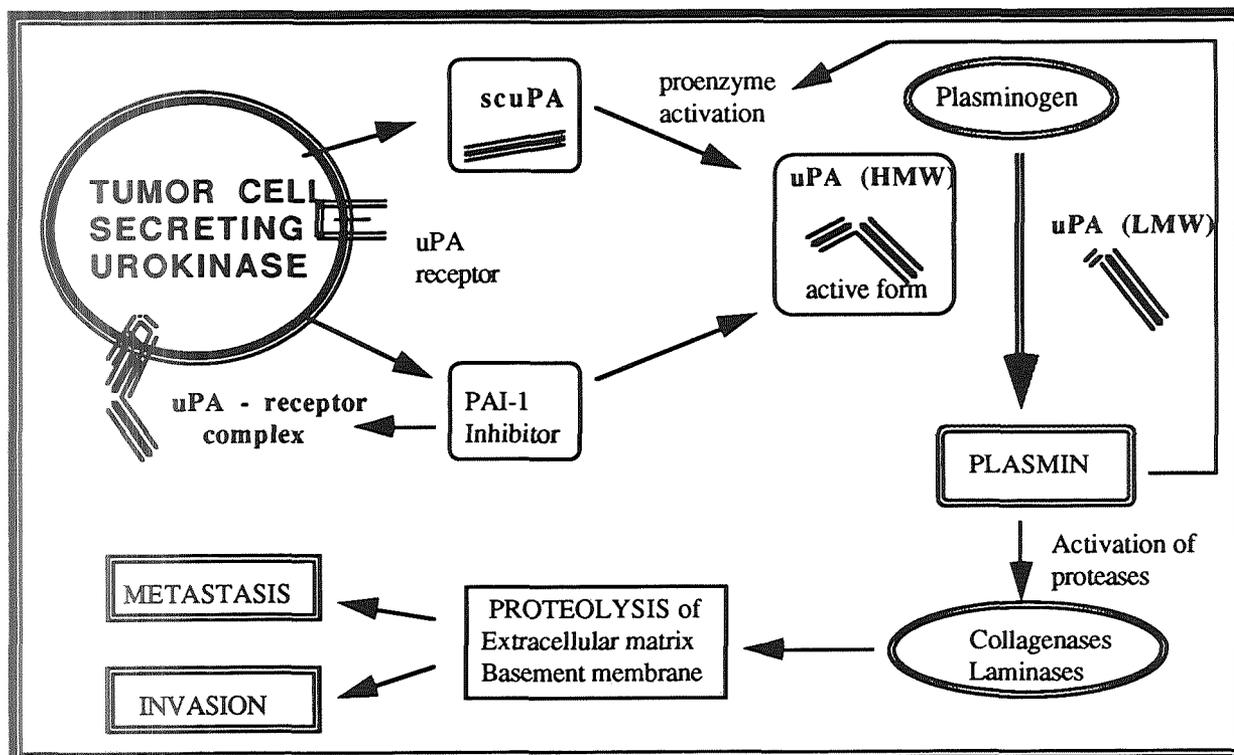


Figure 1. Schematic of interrelationships among urokinase forms, inhibitor, uPA receptor, activation of plasmin and subsequent protease steps that enable tumor cell invasion and lysis of the extracellular matrix and metastasis.

It is clear that the complexity of the many interrelationships within the cascade of proteolytic activations makes it difficult to use an average value for the level of uPA produced by all of the cells in the tumor. Especially, when most of the normal tissue do not produce uPA and many of the tumor cells do not produce uPA unless they are actively metastasizing. The challenge is to quantitatively measure uPA inside and on the surface of the cancer cells and then correlate those uPA levels with other specific markers to characterize the metastatic scenario for each tissue type and stage of cancer. No previous method has been developed to accurately measure the intracellular urokinase content, membrane-bound urokinase and cellular secretion levels and then correlate those urokinase levels with DNA content, DNA synthesis, hormone receptors and other markers of aggressive tumor growth to determine the metastatic potential. This project is developing a quantitative diagnostic test to be used first with existing panels of cytological evaluations of breast cancer and later for other types of cancer.

QUANTITATIVE ANALYSIS OF INDIVIDUAL CELLS

We have used flow cytometry and image analysis of fluorescent microscopic images to measure urokinase and DNA in histopathology tissue sections of breast tumors, dissociated cells (prepared in single cell suspensions) taken from tumor biopsies and in several cell lines of malignant brain tumors (gliomas). Fresh cells are isolated from tumor tissue or cytological samples and prepared for antibody incubation in the same manner. Histology sections are prepared from frozen tissues or deparaffinized sections cut from previously embedded biopsies. The antibodies specific for urokinase are incubated with the cells or tissues first, then the cells are incubated with a second antibody having a fluorescent marker detectable by analytical cytometry techniques. DNA content and synthesis rate (based on DNA stains or uptake of DNA precursors) is measured by flow cytometry or image analysis. The same cell sample can be measured for DNA content and urokinase by staining of the DNA and labeling the urokinase with a fluorescent marker that emits at a different wavelength than the DNA dye or marker. Thus both the DNA and urokinase can be measured simultaneously using two-color image analysis or flow cytometry. The image analysis can localize and quantitate uPA in the cytoplasm and cell membrane. An advantage of the use of cell lines is the ability to study uPA expression in relation to cell proliferation and DNA replication. We also are conducting a retrospective study on biopsies from 500 Stage I, node negative, breast cancer patients in collaboration with the Ontario Oncology Working Group made up of researchers from three Canadian and three U.S. cancer centers.

Methods

Attempts have been made to study the expression of uPA during exponential growth as well as in cultures that have been placed on serum-free medium. Quantitation of uPA levels involves immunofluorescent staining with anti-uPA monoclonal antibody (#394, obtained from American Diagnostica, Greenwich, CN), as primary antibody by the indirect technique. The second antibody consisted of fluorescein-conjugated goat anti-mouse IgG, for FCM and image cytometry studies, or in some cases, rhodamine-labeled goat anti-mouse IgG for image cytometry.

Cells were scraped from flasks, in lieu of trypsinization, in order to preserve membrane-bound antigen. Cells were washed in PBS and triturated to disperse the cell pellets. Single-cell suspensions were usually achieved, which were then fixed for in 0.5% paraformaldehyde 15 min. at room temperature, followed by one hour permeabilization in 70% methanol at 4° C., cells were re-suspended and blocked with 1% bovine serum albumin in PBS for 15 minutes, cells were then washed with PBS and stained for one hour with increasing dilutions of the primary antibody or the equivalent concentrations of naive mouse IgG, both diluted in 1% BSA in PBS. cells were washed and incubated in second antibody diluted in 4% goat serum for one hour. In later experiments, cells were pre-incubated for one hour in 4% goat serum in PBS prior to addition of the second antibody.

Analysis of fluorescence by FCM was conducted with the 488 nm line of an argon laser of a EPICS Profile flow cytometer (Coulter Corporation, Hialeah FL). Green light (from fluorescein emission) was directed by means of a dichroic mirror to pass through a narrow bandpass interference filter (520 +/- 10nm) to impinge on the green-sensitive PMT. Red light (from the DNA stain PI) was deflected through a 630 long pass filter to the corresponding PMT. Bivariate, 64 x 64 channel histograms were obtained for analysis of mean fluorescence intensity.

Digital image analysis was conducted using both Nikon and Zeiss fluorescence microscopes, equipped with a high resolution video camera connected to a QuickCapture board (Data Translation, Inc.) for the MacIntosh II Ci and Fx. The fluorescent filters in the Zeiss microscope were matched closely with the bandpass filters of the EPICS so that image analysis and FCM data on cells from the same sample could be compared. Images were stored as TIFF files and later analysed using NIH Image Version 1.4 (public domain software from NIH). Individual cells were scanned for mean optical densities and normalized for area. Areas of concentrated uPA (including membrane-bound) were further analysed by density slicing and thresholding followed by particulate analysis of those specific areas. Data were also normalized for area and staining intensity after the background was subtracted. This allows comparisons among cells from the same samples and comparisons between cell lines and different samples. Statistical analysis of the data was performed by multivariate analysis using Statview 512 (Abacus Concepts, Inc.)

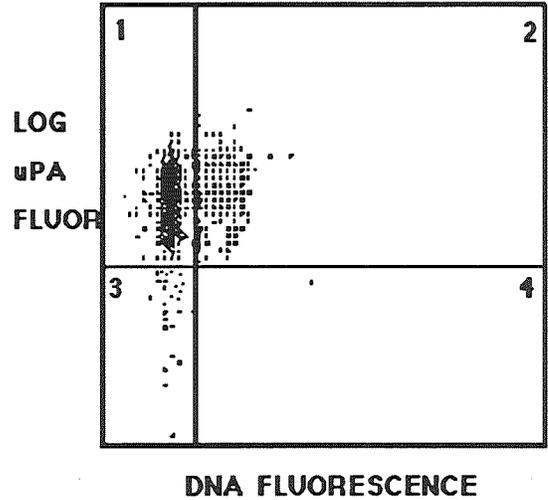
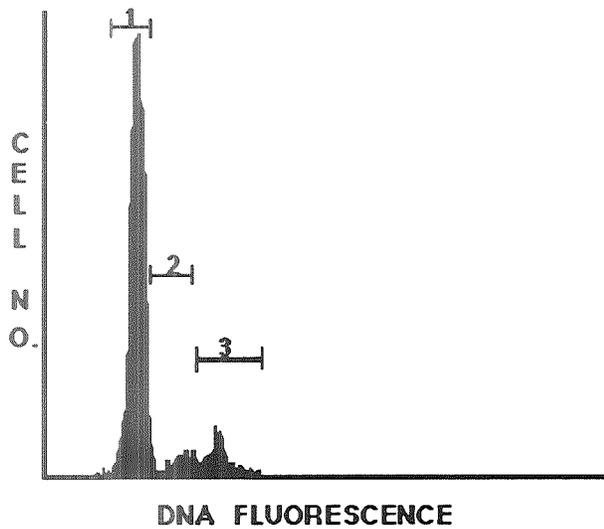
Flow Cytometric Studies of Urokinase in Cultured Glioma Cells:

In order to establish the parameters for immunofluorometric analysis of urokinase (uPA) in tumor cells, studies were initiated with U937 lymphoma and human glioma cell lines, which were found to produce high levels of the plasminogen activator.

The two glioma cell lines employed in the studies were obtained from Dr. Marylou Ingram, Huntington Research Foundation, Pasadena, CA. The two cell lines, which were cultured from patient surgical biopsy material, have different morphological characteristics and growth rates. While alterations in the cells obviously occur in culture, the consistent morphology of these cell lines during passage in culture encouraged us to pursue differences in the cells' characteristics, which can provide correlations between uPA and metastatic relationship of uPA to the biological behavior of the original tumors. The first of these lines, CS, grows very rapidly as polygonal cells in monolayers and, in the absence of serum, tends to form spheroid structures. The second cell line, HBR09, has a fibroblastoid morphology although it has the characteristic immunological marker associated with gliomas, glial fibrillary acidic protein (GFAP) [22]. This cell line grows at about one fourth the rate of CS.

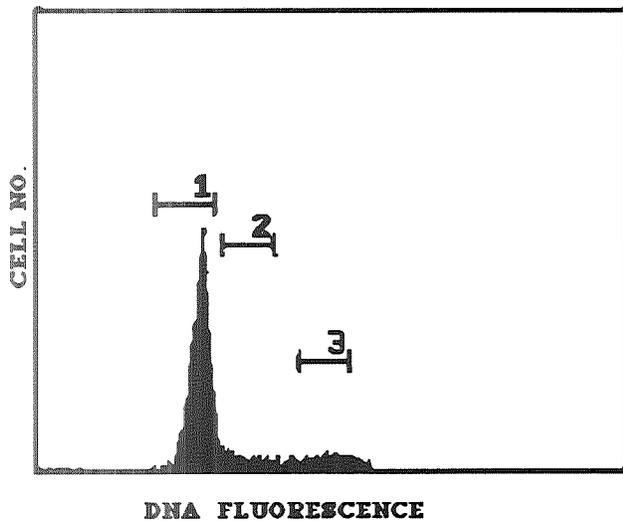
RESULTS

Initially, there were some problems with non-specific fluorescence background which interfered with uPA quantitation. The high background was determined to be due to autofluorescence, since the levels of non-specific fluorescence remained the same even when PBS was substituted for the second antibody. Nevertheless, signal-to-noise was sufficient to measure significant differences in the immunostaining with the anti-uPA Mab. Dual staining with propidium iodide (PI) following ribonuclease treatment and fluorescein-labeled anti-uPA antibodies enabled bivariate analysis of DNA and uPA content as shown in Figure 2.



	MIN.	MAX	PERCENT	MEAN FL.	SD	%HPCY
1 X	0	13	82.6	9.9	1.0	7.26
Y	3.24	1023		12.2	0.17	8.62
2 X	14	63	15.0	17.4	2.2	13.2
Y	3.24	1023		15.48	0.18	10.9
3 X	0	13	2.2	9.5	1.1	9.59
Y	0.102	3.24		12.16	0.26	4.08

Figure 2. Flow cytometry immunofluorescence of glioma cells (HBR09 line) labeled with propidium iodide (PI) for DNA and fluorescein-conjugated antibodies for urokinase (uPA). Panel A shows the DNA histogram of these cells with G₁, S and G₂ + M subpopulations. Panel B compares the uPA and DNA fluorescence for cells in G₁ (82% of total), in S phase (15% of total) and in G₂ + M (22% of total).



#1141
CS uPA

LOG uPA FLUORESCENCE

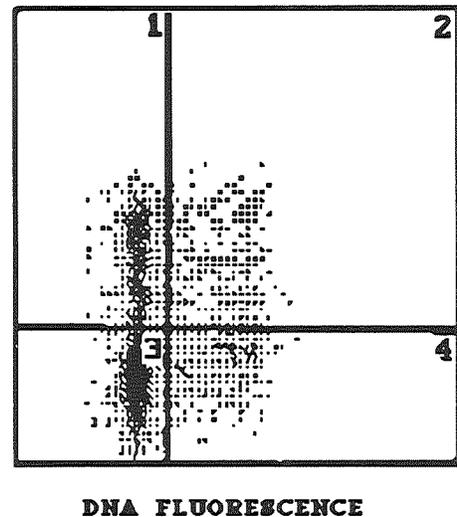


Figure 3. Flow cytometry analysis of CS cell line for comparison with Figure 1 above. Panel A shows the CS cell cycle distribution of DNA fluorescence (PI). Panel B shows the fluorescence distribution of uPA vs. DNA. Most of the cells in G₁ phase (22% of total) and in S phase (7% of total) contain significant levels of uPA.

We are comparing flow cytometric analysis of uPA levels in both CS and HBRO9 glioma lines. The relative levels of uPA as measured by flow cytometry immunofluorescence are tabulated in Table 2. Further studies are concentrating on measurements of uPA by image cytometry, in an effort to distinguish the membrane-bound (receptor) vs cytoplasmic uPA, since flow cytometry only measures fluorescence at "zero resolution". FCM studies have demonstrated that two color fluorescence can be used to measure uPA and DNA in the same cell population. There often are some non-standard cells in some individual cultures that require careful placements of the gates to get representative cells in all three phases of the cell cycle. It is important to know that the measurement methods are sensitive enough to determine the variability among replicate cultures from the same tumor source, since there is always some degree variability from patient to patient in any marker expression.

Both glioma cell lines produce uPA during growth and also during stationary (G_1) phase, HBRO9 producing significantly higher levels of intracellular uPA. They also appear to produce more membrane-bound uPA (based on qualitative examinations of some 150 cells), however, quantitative measurements are still underway. It is noted that the HBRO9 cells had considerably more variability in the fluorescence measurements than did the CS cells.

Table 2. Urokinase levels in glioma cell lines measured by immunofluorescence flow cytometry.

<u>CELL LINE</u>	<u>MEAN FLUORESCENCE</u>	<u>+/- s.d.</u>
CS	6.13	0.10
HBRO9	32.25	11.3

Image Analysis of Fluorescent-labeled uPA in Breast and Brain Tumors

Evaluations of anti-uPA labeled breast cancer sections reveal that normal breast tissue does not contain uPA except for some endothelial cells lining the arterioles. Intraductal carcinomas, however, do express measurable quantities of uPA [23]. Quantitative measurements of uPA by absorption of immunologic stains as light passes through tumor cells is difficult since histopathology counterstains add to the absorption of the uPA antibody labels. Fluorescence is a better quantitative tool since the light is emitted only from the uPA molecules and it is emitted at a wavelength different from the incident light. Fluorescence emitted from whole cells (cytoprep) can clearly show the concentrations of uPA on the cell membrane as well as "hot spots" within the cytoplasm. Figure 4a shows an example of a breast tumor section illustrating the areas of uPA found in foci of tumor cells. Distinct areas of concentrated uPA are shown (white lines). Clearly, many tumor cells are not producing significant quantities of uPA and neither are most of the normal cells. Thresholding and image enhancements can often give size distribution and more information on the cells producing the uPA. Figure 4b shows the same tissue section as Figure 4a, however, this image has been analysed and pseudocolor added to the display to illustrate that considerable cellular detail remains obscured in the photo 4a.

Glioma cells that have been labeled with rhodamine-conjugated antibodies for uPA are shown in Figure 5a. Note that these cells (CS line) contain large amounts of urokinase per cell and that the uPA concentration is quite varied throughout the cytoplasm and there is a lot of cell to cell variation. Each cell can be scanned to obtain optical density levels that can be compared among cells after normalizing with area of each cell measured. Some cells exhibit "hot spots" of concentrated uPA, especially on the membrane. It is possible to measure the membrane-bound uPA by differential analysis using the mean density level of the weaker cytoplasm subtracted from that of the whole cell containing the membrane bound enzyme. Specific areas of uPA also can be measured by selecting a density slice(s) that represent the major portion of concentrated uPA. The particle size (number of pixels) to be counted is defined, then that particular density slice of fluorescence can be measured automatically by particle analysis. This will give the number of particles, average particles per group, area, perimeter of the cells and location of particle groups larger than a defined size (see Figure 5b).

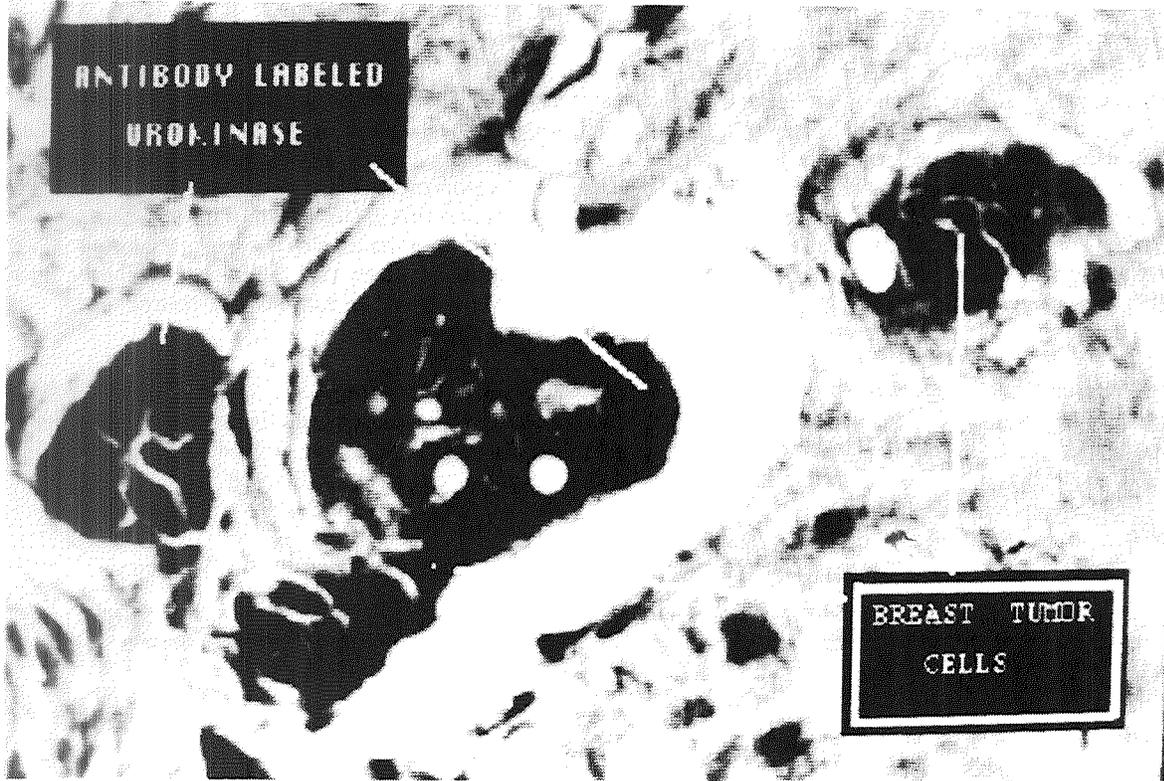


Figure 4a. Digital image of breast histology section showing tumor cells and antibody labeled urokinase areas (white lines) selected by density slicing for quantitative measurement.

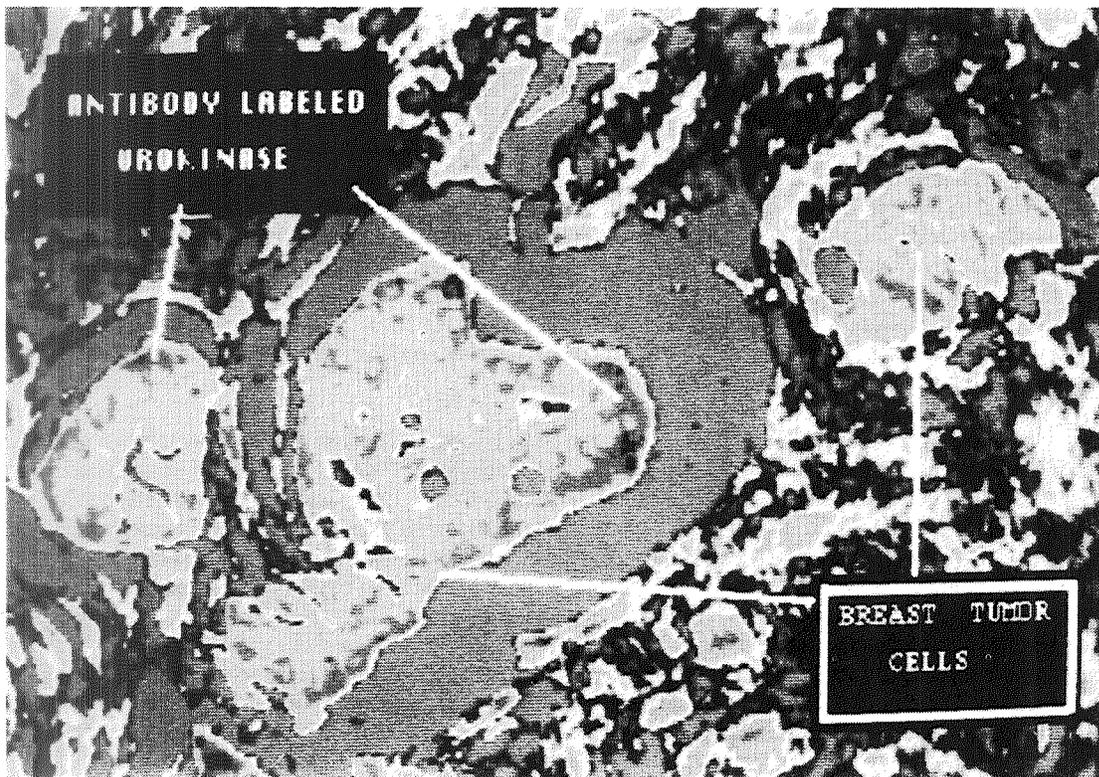


Figure 4b. Pseudocolor image of same section as 4a above, showing the additional morphological information contained in the recorded grey levels in the original digitized image.



Figure 5a. Photomicrograph of human glioma cells stained for urokinase (uPA) by rhodamine labeled antibodies. Note areas of concentrated uPA in selected areas of the cytoplasm and in some areas of the cell membranes.

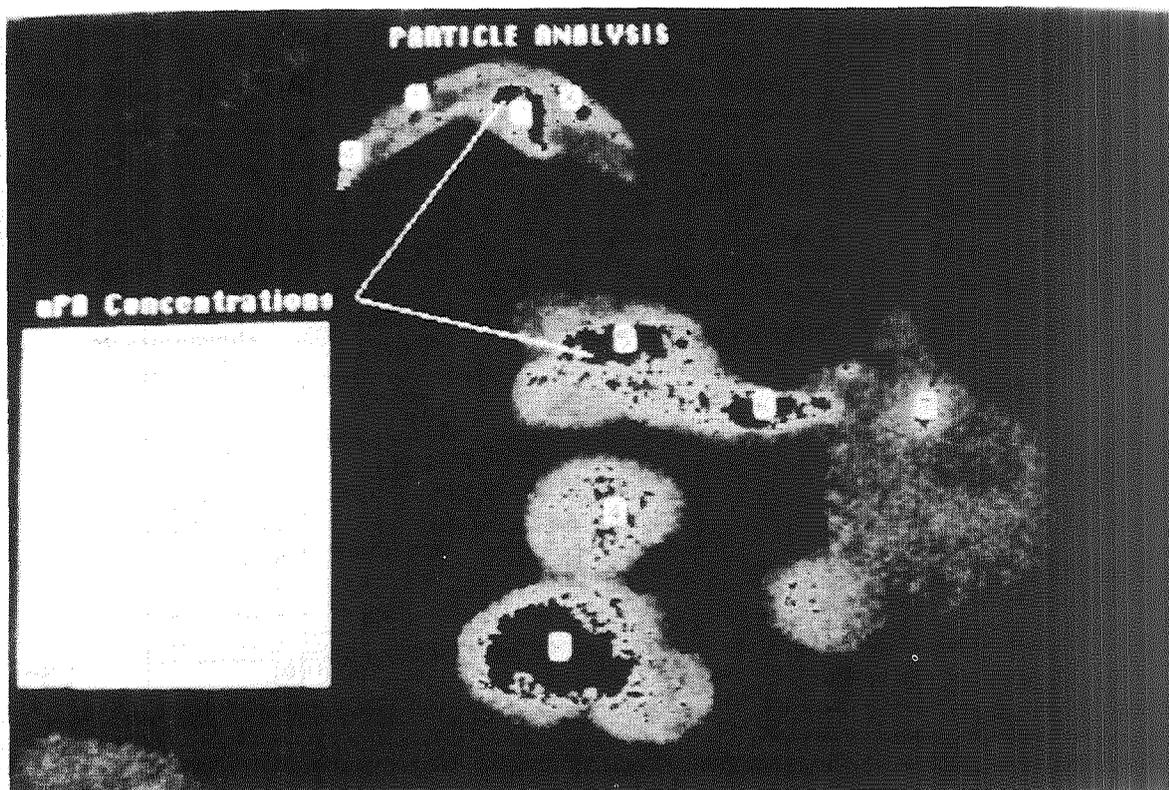


Figure 5b. Digital image of glioma cells labeled with rhodamine-conjugated anti-urokinase antibodies. Image has been density sliced, particle size selected from 50 to 5000 pixels, then particle analysis performed, with the major areas of urokinase counted and labeled. The area and mean optical density (related to fluorescence intensity) is also recorded for statistical comparisons between cell lines and patient biopsies.

DNA can also be quantitated on a per-cell basis using image analysis. However, when using PI for DNA and fluorescein label for uPA, the amount of DNA fluorescence was often more predominant than was the uPA-related fluorescence (green wavelength) since PI fluorescence overrides FITC emissions. This required adjustment of the incident light intensity to keep both fluorescence signals in the same range so as to avoid resetting the video camera sensitivity between measurements on the same field of view. DNA quantitation can be performed more effectively by staining a dye excited in the U.V. (Hoechst 33258) and analysed in the blue region.

DISCUSSION

The importance of urokinase as a key enzyme in the initial mechanisms leading to tumor cell invasion and metastasis has been underscored in the past three years. Previous methods of measuring extracted uPA /mg. of protein or measuring secretion levels in cultured explants have provided statistical correlation with disease-free and overall survival [23]. There also is a strong correlation between uPA production and lymph node status in breast cancers and multivariate analyses have shown that high levels of both uPA and PAI-1 means a maximum risk of relapse. It is now time to develop more specific tests that can accurately determine the active uPA vs. the inactive scuPA, the membrane bound uPA and the PAIs that appear to have interlinked, critical roles in the migration and metastasis of breast and other cancers.

Correlations of uPA with other markers require more precise knowledge about uPA and the multiple biochemical interactions that affect its proteolytic action. Figure 6 illustrates the current methods of measuring average levels from all tumor cells vs. our method for measuring uPA directly in the cells. These methods can be used in retrospective studies where the time to reoccurrence, degree of metastasis and morbidity are known. Cumulative data on many patients (>50) can then be used to provide a prognostic indicator for the presence and degree of active metastasis occurring in primary tumors. A study of uPA in node negative breast cancer is underway.

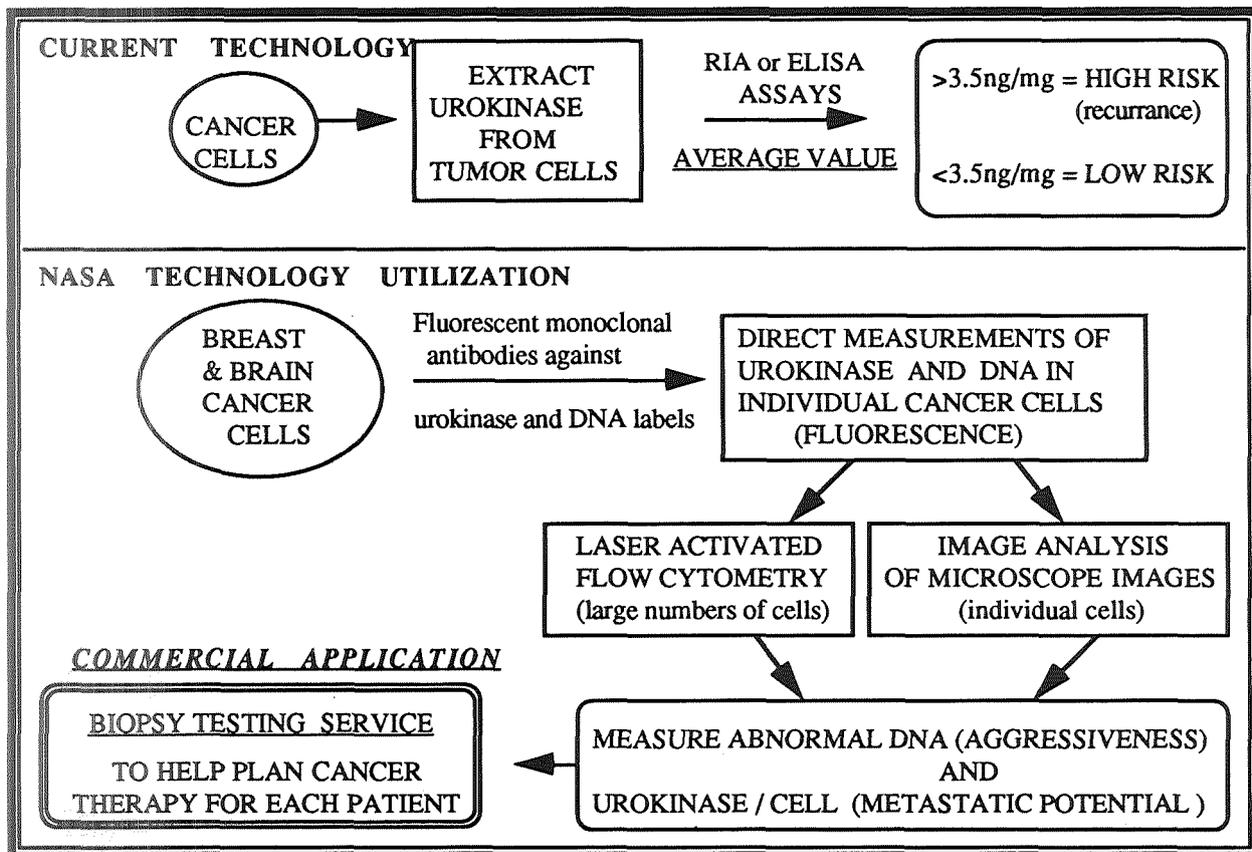


Figure 6. Schematic of new methods for quantitative measure of tumor markers and metastatic potential.

Research validation: correlations of uPA and other markers

The first research step has been to compare the DNA measurements and the intracellular levels of urokinase in tumor cells and normal cells. The initial FCM analyses will determine the effect of cell cycle on those cells having elevated uPA and the general relationship between abnormal DNA and uPA in breast and brain tumors. Additional data is being collected on DNA synthesis rates and uPA levels using more specific flow cytometry and/or image analysis techniques. Urokinase levels can be determined in those subpopulations of tumor cells that have abnormal DNA (using two parameter flow cytometry). We currently are measuring the intracellular and membrane-bound levels of urokinase per cell using fluorescent anti-uPA antibodies and image analysis. Next, different forms of uPA and uPA receptor complexes will be measured using molecular-specific antibodies. Finally, PAI-1 or PAI-2 will be measured per cell and compared to the abnormal DNA and high uPA to determine the final relative metastatic potential. Correlations with hormone receptors and other proteolytic enzymes also can be made to provide additional prognostic information for custom design of adjuvant therapy following surgical removal of the primary tumor.

COMMERCIAL APPLICATIONS

Many intricate biochemical interactions are involved in dissolution of the extracellular matrix which enables metastatic cells to leave the primary tumor. Intracellular metabolism appears to have a major role in the initiation of cellular metastasis. The complexity of these interactions makes it too complicated for laboratory test kits to be effective in routinely measuring the DNA, uPA, PAIs, hormone receptors, etc. necessary to develop a comprehensive prognostic panel for a particular cancer patient. Such a task requires a specially equipped, expert medical testing service where pathologists, surgeons and oncologists can send patient biopsies for complete analysis. Several companies already offer this type of cancer testing as a commercial service, however, tests for uPA as a prognostic marker of metastatic potential are not offered yet. Once the metastatic relationships are characterized at the cellular level, clinical studies will be required to statistically correlate these biochemical tests with recurrent disease and survival. Quantitative measurements of uPAs, uPA receptors and inhibitors can be added to the existing panel of breast cancer cytological tests. The first use of these tests will be in providing additional information that indicate active metastasis at the time of initial surgery. This information can help oncologists design better, more effective follow-up therapy for those patients that have high levels of multiple markers indicating metastasis is already underway even though clinical manifestations are still undetected.

This NASA sponsored project is developing methods for a routine analytical test of intracellular and membrane-bound uPA that can be added to the existing panel of breast cancer markers. Each year more than 170,000 new breast cancers are discovered in the U.S. alone. Unfortunately, about 30% of these patients will die from their breast cancer [24]. The current tests for DNA content, DNA synthesis and hormone receptors cost about \$350. More complicated tests will likely cost \$450 each. A practical test for urokinase combined with other metastatic markers of breast cancer would create a significant new market for cancer testing laboratories. And of course, uPA is important in many other types of metastatic cancers. Better adjuvant therapy, used only when critical markers are known to indicate active metastasis, could make a significant impact on the survival of cancer patients and reduce medical costs required to treat recurrent disease.

Finally, the use of antibodies specific against urokinase can be used for more than diagnosis of the beginning steps of metastasis. As the entire scenario is better understood, it may be possible to develop treatments targeted against just those metastatic cells that have large amounts of membrane-bound urokinase or large concentrations of inactive scuPA. Anti-uPA antibodies could be conjugated with anti-tumor drugs or radioisotopes to treat specific metastatic cells that are actively trying to invade adjacent tissues. This could be the basis for the first therapy directed against metastatic cells that were not removed by cancer surgery or began migration prior to removal of the primary tumor.

REFERENCES

1. Gratzner, H.G. and Leif, R.C., "An immunofluorescence method for monitoring DNA synthesis by flow cytometry," Cytometry 1: 385-389, 1981.
2. Merkel, D.E. Dressler, L.G. and McGuire, W.L., "Flow cytometry, cellular DAN content and prognosis in human malignancy," J. Clin. Oncology 5:1690-1703, 1987.
3. Celis, J.E. and Celis, A., "Cell cycle-dependent variations in the distribution of the nuclear protein cyclin proliferating cell nuclear antigen in cultured cells: subdivision of S phase," Proc. Natl. Acad. Sci. 82: 3262-3266, 1985.
4. Lee, J.S. et al., "Determination of biomarkers for intermediate end points in chemoprevention," Cancer Research, Vol. 52 (Suppl.): 2707s - 2710s, 1992.
5. McQuire, W. L. "Prognostic factors in primary breast cancer," Cancer Surveys 5: 527-536, 1986.
6. Raber, M.N. et al., "Ploidy, proliferative activity and estrogen receptor activity in human breast cancer," Cytometry, 3: 36-41, 1982.
7. Bacus, S. S. et al., "Tumor-inhibitory monoclonal antibodies to the HER-2/neu receptor induce differentiation of human breast cancer cells." Cancer Research 52: 2580-2589, 1992.
8. Duffy, M.J., et al., "Urokinase-plasminogen activator, a marker for aggressive breast carcinomas," Cancer 62: 531-533, 1991.
9. Janicke, F. et al., "The urokinase-type plasminogen activator (u-PA) is a potent predictor of early relapse in breast cancer," Fibrinolysis 4: 69-78, 1990.
10. Oka, T. Ishida, T. Nishino, O. T. and Sugimachi, K., "Immunohistochemical evidence of Urokinase-type plasminogen activator in primary and metastatic tumors of pulmonary adenocarcinoma," Cancer Research 51: 3522-25, 1991.
11. Hollas, W., Blasi, F. and Boyd, D. "Role of urokinase receptor in facilitating extracellular matrix invasion by cultured colon cancer," Cancer Research 51: 3690-3695, 1991.
12. Harvey, S.R., et al. "Secretion of plasminogen activators by human colorectal and gastric tumor explants." Clin. Expl. Metastasis 6: 431-450, 1988.
13. Camiolo, S.M., Markus, G. and Piver, M.S. "Plasminogen activator content of gynecological tumors and their metastases." Gynecological Oncology 26: 364-373, 1987.
14. Hasui, Y. et al. "Comparative study of plasminogen activators in cancers and normal mucosa of human urinary bladder." Cancer Research 49: 1067-1070, 1989.
15. Camiolo, S.M. et al., "Plasminogen activator content and secretion in explants of neoplastic and benign human prostate tissues," Cancer Research 44: 311-318, 1984.
16. Meissauer, A. et al., "Urokinase-type and tissue-type plasminogen activators are essential for in vitro invasion of human melanoma cells," Experimental Cell Research 192: 453-459, 1991.
17. Schmidt, M. et al. "Tumour-associated fibrinolysis: the prognostic relevance of plasminogen activators uPA and tPA in human breast cancer." Blood Coagulation and Fibrinolysis 1: 695-702, 1990.
18. Gaylis, F.D. et al., "Plasminogen activators in human prostate cancer cell lines and tumors: Correlations with the aggressive phenotype," Journal of Urology 142: 193-198, 1989.
19. Kirchheimer, J.C. and Remold, H.G., "Functional characteristics of receptor-bound urokinase on human monocytes: Catalytic efficiency and susceptibility to inactivation by plasminogen activator inhibitors," Blood 74: 1396-1402, 1989.
20. Andreasen, P.A. et al., "Plasminogen activator inhibitor from human fibrosarcoma cells binds urokinase-type plasminogen activator, but not its proenzyme," J. Biol. Chem. 261: 7644-7651, 1986.
21. Cubellis, M.V., Wun, T. and Blasi, F., "Receptor-mediated internalization and degradation of urokinase is caused by its specific inhibitor PAI-1," EMBO Journal, 9: 1079-1085, 1990.
22. Patau, A., "Glial GFAP, vimentin and fibronectin in primary cultures of human glioma and fetal brain," Acta Neuropathology 75: 448-455, 1988.
23. Katz, R., M.D. Anderson Cancer Center, personal communication, Dec. 1991.
24. Janicke, F., Schmitt, M. and Graeff, H. "Clinical Relevance of the urokinase-type and tissue-type plasminogen activators and of their type 1 inhibitor in breast cancer," Seminars in Thrombosis and Hemostasis 17: 303-312, 1991.
25. Fisher, B. "The evolution of paradigms for the management of breast cancer: a personal perspective," Cancer Research 52: 2371-2383, 1992.

IMMUNOCONJUGATES: "MAGIC BULLETS" FOR CANCER THERAPY?

Daniel R. Passeri and Jack Spiegel
Office of Technology Transfer
National Institutes of Health
Bethesda, MD 20892

N 9 3 - 5 2 5 3 6 7
150476
p. 6

Abstract

Conjugating cytotoxic agents to antibodies allows for site-specific delivery of the agent to tumor cells and should provide increased efficacy and reduced non-specific toxicity. These site-specific cytotoxic agents are known as immunoconjugates or "magic bullets" and have demonstrated great promise as therapeutic agents for cancer and other diseases. The historical developments and future potential of this new approach to cancer therapy are reviewed.

Introduction

Cancer is responsible for approximately twenty-five percent of the deaths in industrialized countries and it is estimated that there are currently over seven million cancer patients in America. The American Cancer Society estimates that over one million Americans will be diagnosed with cancer in 1992 and that approximately 520,000 people will die of cancer, making it the second leading cause of death in America. Despite these alarming statistics, anticancer therapeutics account for only the eighth-largest pharmaceutical market in the United States.

Cancer causes far more morbidity and mortality than diseases that account for far larger drug markets and has shown much slower market growth. The slow market growth for the cancer therapeutic market is primarily because of the problems associated with low efficacy and the serious side effects of the majority of anticancer drugs. Even the most effective drugs, including some biotherapies, show as little as a fifty percent success rate. In addition, most cancer therapy is extremely toxic since chemotherapeutic agents kill normal cells as well as cancer cells.

Cancer therapies have historically consisted primarily of surgery, chemotherapy, and radiation. Because of the risks and invasive nature of surgery, and the adverse effects of radiation and chemotherapy, there is tremendous opportunity for new non-invasive therapies which offer improved efficacy while reducing associated side effects. Recent advances have provided opportunities for developing new alternative treatment strategies. One approach is to target a cytotoxic agent to the cancer cell through the development of immunoconjugates. The term immunoconjugate designates monoclonal antibodies (or antibody fragments containing their binding sites) linked to cytotoxic agents: drugs, toxins, radioisotopes or cytotoxic cells of the immune system. To accomplish this, the cytotoxic agent is attached to an antibody or a growth factor that preferentially binds to cancer cells.

This exciting new technology has proven to be a feasible alternative to conventional chemotherapy and shows promise as an effective therapy for many cancers which have not responded well to conventional therapies. According to a 1991 Frost & Sullivan market report, the market for cancer therapy immunoconjugates is projected to be approximately \$720 million by 1996 and should continue to grow as technological improvements provide for higher clinical efficacy and general market acceptance [1].

This article will review the progress as well as several problems inherent to this approach to cancer therapy, and briefly highlight approaches taken at the NIH and elsewhere to advance this technology.

Historical Background and Current Developments

About 1913, Paul Ehrlich conceived the idea of therapeutics which function as "magic bullets" [2]. Ehrlich's "magic bullets" combined the targeting properties of antibodies with cytotoxic agents. Nearly eighty years later, this vision is beginning to reach fruition. The protracted period between Ehrlich's conception and the potential realization of immunoconjugate therapy for certain cancers underscores the numerous technological difficulties encountered in this field.

The first consideration in developing a "magic bullet" for cancer therapy is selection of an appropriate carrier or targeting agent to deliver toxic agents specifically to the tumor cells. In theory, the targeting agent could be any moiety capable of selective binding to a receptor on tumor cells. Indeed, anti-tumor reagents have been produced by attaching cytotoxic agents to numerous cell recognition proteins, including antibodies, alpha transforming growth factor, epidermal growth factor, interleukins, and transferrin [3,4]. Site-directed cytotoxicity is aimed primarily at cell-surface antigens or at receptors expressed in high numbers on cancer cells or other cells of interest. Toxic substances can be conjugated to antibodies or fusion proteins that recognize the cell-surface antigens characteristic of the specific cell type that is targeted for treatment. The toxin complex then specifically binds to the targeted cells resulting in a localized high dose of the toxin to them while sparing the normal cells. Early attempts to develop "magic bullets" using polyclonal antisera against tumor cells were frustrated by significant cross reactivity with surface antigens on normal cells. Even after rigorous absorption against normal tissues, polyclonal antisera preparations vary markedly in reactivity, specificity, and reducibility [5]. Coupling cytotoxic agents to such antibodies exacerbates the non-specificity, and results in therapeutic preparations with unacceptable toxic side effects.

Development of antibodies specific to tumor cells was revitalized with the introduction of hybridoma technology in 1975 which permits production of monoclonal antibodies against a selected antigen [6]. Despite this significant advance over polyclonal antisera, many monoclonal antibodies to tumor cells retain some degree of cross-reactivity with normal cells. Anti-tumor immunotoxic conjugates displaying cross-reactivity to normal tissue must be employed judiciously to minimize toxic side effects. The search for new monoclonals with greater specificity against tumor cells is an ongoing endeavor. For example, NIH scientists have patented or have pending patent applications claiming an expanding portfolio of selective monoclonal antibodies useful for treating a wide variety of cancers, including medulloblastoma, glioblastoma, adenocarcinomas, squamous cell carcinomas, breast, colon, prostatic, ovarian, cervical, and esophageal cancer.

Most currently available monoclonals against tumor cells are derived from murine hybridomas. Multiple administration of such monoclonals stimulates immunological responses by the human host against the foreign mouse immunoglobulins [7]. Neutralizing human anti-mouse antibodies compromise the efficacy of immunoconjugate therapy. The ultimate solution to this problem, of course, would be to utilize human monoclonal antibodies rather than murine species. Progress continues in the developing field of human hybridoma technology. A group from the National Cancer Institute and Bionetics Research (a division of Organon-Teknika) has reported testing the immunogenicity of conjugates incorporating two human immunoglobulins directed against colorectal cancer. To date, the study has confirmed the general expectation that immunogenicity of human antibodies will be low [1]. However, nagging technical problems remain with establishing human hybridomas which prevent them from being a reliable source of human monoclonals in the near term.

In lieu of human hybridomas, much work has been directed toward "humanizing" murine monoclonal antibodies. A simple technique which markedly reduces anti-mouse immunoglobulin effects is to use only those portions of the immunoglobulin molecule responsible for binding affinity and specificity. Consequently, immunoconjugates have been constructed using Fab, Fab' or F(ab')₂ fragments rather than intact immunoglobulin [8]. This eliminates the immunogenic epitopes of the Fc region of the mouse antibody. Two additional benefits may accrue from removal of the Fc portion of the antibody. Firstly, large immunoglobulin conjugates have difficulty permeating solid tumors. This problem is reduced substantially when smaller antibody fragments are employed. Secondly, such fragments eliminate non-specific binding to

non-target cells mediated via the Fc region; e.g., binding to cells of the reticuloendothelial system. These potential benefits must be weighed against negative consequences of using antibody fragments. For example, antibody fragments are known to be cleared from circulation more quickly than intact immunoglobulins [9]. Additionally, Fab fragments would not be indicated in situations where the Fc portion of the antibody is critical to the biological function of the antibody. For example, Michael Kaliner of the National Institute of Allergy and Infectious Diseases has developed a model for selective destruction of cells expressing high affinity IgE Fc receptors (e.g., mast cells in either malignant systemic or benign systemic mastocytosis) employing IgE immunotoxin conjugates [10].

Another approach to "humanize" murine monoclonals for human therapy relies upon application of recombinant DNA technology. Chimeric (or mouse/human) antibodies have been created whereby the constant regions of human immunoglobulins are fused to the variable regions of mouse monoclonal antibodies [11]. These chimeric antibodies retain the antigen binding specificity of the mouse monoclonal, elicit reduced human anti-mouse antibody responses in patients, and are not subject to the enhanced clearance rates of Fab fragments.

An exquisite extension of the recombinant approach involves constructing mouse/human chimeric antibodies which incorporate only the complementarity-determining regions (CDRs) from the mouse. CDRs are the portions of the antibody molecule which guide the antibody to its binding ligand. The remainder of the chimeric antibody structure is human, including the framework residues (FR) which support the CDRs and determine the disposition of the CDRs relative to one another [12]. A further variation on this technique, called "veneering", was developed as a joint invention by Merck & Co., Inc. and Eduardo Padlan of the National Institute of Diabetes and Digestive and Kidney Diseases. Veneering judiciously replaces mouse exterior amino acid residues in the variable region of the antibody with those of the human. The premise of "veneering" is that the key residues in CDRs (i.e., those involved in preserving ligand binding) are "interior" and interdomain contact residues. Consequently, surface amino acid residues of mouse origin, which can be "seen" by the immune system in its immune surveillance, may be changed to their human counterpart without affecting ligand binding properties.

Selection of appropriate and optimal cytotoxic components for immunoconjugates also has been an area of active research. Early anti-tumor immunoconjugates combined antibodies with known low molecular weight chemotherapeutics such as radionuclides, DNA alkylating agents, and anti-metabolites. Recently, Otto Gansow of the National Cancer Institute has reported encouraging results using yttrium-90 conjugated to anti-interleukin-2 receptor antibody for treating T-cell leukemia patients. The greatest interest, however, has been in the use of bacterial and plant toxins as the cytotoxic element of therapeutic immunoconjugates. The best studied of these toxins are diphtheria toxin (DT) from *Corynebacterium diphtheria*, the lectin ricin from the seeds of *Ricinus communis*, and pseudomonas exotoxin A (PE) from *Pseudomonas aeruginosa*.

Both DT and ricin are heterodimeric molecules consisting of A- and B-chains. In both toxins, the B-chain is responsible for cellular binding and entry into the target cell, and the A-chain is a potent inhibitor of protein synthesis. DT bound to cell surfaces by the B-chain enters the cell via receptor-mediated endocytosis. Within the resulting endosomes, the B-chain of diphtheria toxin undergoes a conformational change which permits the A-chain to translocate into the cytoplasm. Once in the cytoplasm, diphtheria toxin A-chain irreversibly inhibits the protein translation machinery. Specifically, A-chain inactivates elongation factor 2 (EF-2) via an ADP-ribosylation reaction. Ricin binds to cells via affinity of its B-chain for terminal galactose residues of cell surface glycoproteins and glycolipids. Analogous to diphtheria toxin, surface-bound ricin becomes internalized within vesicle structures, and the B-chain facilitates the translocation of the ricin A-chain ("ricin A") out of vesicles into the cytoplasm. Ricin A inhibits protein synthesis via selective N-glycosidase activity which cleaves a specific adenine residue in the 28S ribosomal subunit. Both toxins are extremely potent: a single molecule is sufficient to inhibit protein synthesis within a cell.

Both DT and ricin have been conjugated to antibodies producing immunotoxins with potent cytotoxic activities [13]. Such immunotoxins, however, exhibit serious non-selective binding due to the binding

properties of their respective B-chains. Attempts to remedy this problem by conjugating only the A-chain of the toxin to the antibody produces immunotoxins with good selectivity, but variable cytotoxic potency [14]. Genes for both diphtheria toxin and ricin have been cloned, and recombinant constructs containing mutant B-chains are being tested for reduced cell binding. Richard Youle and colleagues at the National Institute of Neurological Disorders and Stroke have developed a recombinant DT with reduced cell binding properties [15]. When conjugated to anti-human transferrin receptor or anti-CD3 antibodies, this recombinant DT demonstrated up to 1,000 fold reduction in cell binding; yet was equal to wild-type immunotoxin in cytotoxic potency. David Neville's research group at the National Institute of Mental Health has developed similar immunotoxins utilizing recombinant DT.

Another immunotoxin system demonstrating exceptional promise has been developed in Ira Pastan's laboratory at the National Cancer Institute. This system utilizes the bacterial toxin pseudomonas exotoxin A (PE). Analogous to the B and A-chains of DT and ricin, PE contains domain I and domain III which confer cell binding and cytotoxicity, respectively [16]. Cytotoxicity is accomplished by the same mechanism as in diphtheria toxin; i.e., ADP-ribosylation of elongation factor 2. Unlike ricin and DT, pseudomonas exotoxin A additionally has a domain II which mediates translocation of the toxic domain III across cell membranes. Consequently, it has been possible to abolish cell binding (i.e., domain I function) without disturbing membrane translocation functions [17]. Recombinant pseudomonas exotoxins, with truncations in domain I (PE40 and PE38), have been prepared. These modified forms of pseudomonas exotoxin are as potent as native PE, but are 100 fold less toxic to nontarget cells [3,18].

PE40 and PE38 have been used also to construct completely recombinant immunotoxins by fusing them to DNA fragments encoding growth factors, antibodies, and antibody-fragments. In this way, PE40 and PE38 have been joined to the carboxyl end of the fragment variable (Fv) portion of antibodies to produce, so called, "recombinant single chain immunotoxins". The Fv region is the smallest antibody fragment capable of binding antigen. They consist of two chains, each about 110 amino acids in size, held together by a linking peptide about 15 amino acids in length. Recombinant single chain immunotoxins have been constructed to selectively bind the human transferrin receptor, human IL-2 receptor, and an antigen, recognized by monoclonal antibody B3, found on many human carcinomas (e.g., prostate, colon, stomach, breast, ovary, lung, and bladder). These recombinant immunotoxins are particularly attractive in that they can be produced in large amounts in *E. coli*, have reduced animal toxicity, and appear to be well suited to penetrate solid tumors by virtue of their small size [19-21].

The application of recombinant toxins has markedly reduced toxic side effects associated with the native molecules. These still are foreign proteins, however, and repeated administration leads to host immunological responses against the toxin. Richard Youle's laboratory has developed an approach to further "humanize" immunotoxins. They constructed a recombinant immunotoxin, where the traditional bacterial or plant toxin is replaced by a human enzyme, angiogenin [22]. Angiogenin is a protein found in normal blood plasma, and has homology to pancreatic Rnase. While not cytotoxic toward intact cells, angiogenin is a potent inhibitor of protein synthesis once it gains access to the protein synthesis machinery within the cytoplasm. Attachment of angiogenin to antibodies directed against cell surface antigens results in endocytotic incorporation, followed by inhibition of protein synthesis. Using recombinant techniques, angiogenin was fused to a mouse/human chimeric antibody heavy chain gene. This antibody-angiogenin fusion protein was introduced into a transfectoma which secreted the chimeric light chain of the same antibody [23]. The resultant F(ab')₂-like antibody-angiogenin fusion protein has the "magic bullet" properties of antibodies linked to plant/bacterial toxins, but elicits a reduced immune response in the host.

Also, recent studies are beginning to demonstrate synergistic antitumor effects when immunoconjugates are used in conjunction with other treatment modalities. Because of the advantages of site directed specificity and the potential for synergistic effect, immunoconjugates are expected to replace or supplement, in increasing measure, the use of unconjugated chemotherapeutics and radionuclides in therapy [1].

Market Outlook

The primary criteria for an immunotoxin are specificity and high potency, i.e, the toxin must be delivered to a specific cell type and must be able to get into the cells for maximum cytotoxic effect. The immunotoxin market is currently dominated by the use of one of three toxins, i.e., ricin toxin, pseudomonas exotoxin, and diphtheria toxin.

Although immunoconjugates show much promise for cancer therapy, they will certainly not replace surgery as the primary therapy whenever surgery is feasible. Even when it is known that not all malignant sites can be resected, it is important to decrease the tumor burden as much as possible, so that the non-invasive treatments, and the patient's own defenses, have a reduced task. Large tumor masses, which are most amenable to surgical resection, are also least accessible to conjugates and other pharmacological agents, because of their poor internal circulation [1]. Immunoconjugate therapy will most likely be the primary therapy of choice for inoperable cancers and for cancer micro-metastases, i.e., when cancerous cells move to various locations throughout the body.

Currently, the major obstacle to overall market acceptance of immunotoxins for therapeutic applications is non-specific toxicity. Price is also an obstacle: initially the cost will be \$2,000 - \$5,000/course of treatment for therapeutic immunoconjugates [1]. This increased cost for treatment may however, be justified by the decrease in required hospital care, i.e., the immunoconjugate therapies may prove to be more effective and efficient for treating numerous diseases with fewer side effects, and consequent faster discharge from the hospital.

What role immunoconjugates will play in future cancer therapies is not yet clear, but in solid tumors therapeutic immunoconjugates will probably be most useful for eliminating residual or occult sites of malignancy after surgery, and to reduce the tumor burden, prolong life and improve the quality of life for patients with advanced disease [1].

Immunoconjugate therapy is most obviously applicable to cancer therapy, but may have significant market penetration in numerous other therapeutic applications including: rheumatoid arthritis, diabetes, infectious diseases, AIDS, and graft vs. host disease.

Scientists at the NIH have played a major role in the development of this promising technology and continue to be at the forefront of new discoveries in the fight against cancer. These new developments are brought forward to the market place through the patenting and licensing efforts of the Office of Technology Transfer at NIH. Technology transfer is the process by which the discoveries of laboratories are brought forth into practical knowledge and useful products. The NIH Office of Technology Transfer's primary mission is to facilitate the transfer of technology from Federal laboratories into the private sector for further development and commercialization for the benefit of world health.

It is critical to the medical community and the public welfare that these new technologies find their way to the market place as quickly and safely as possible. The continued efforts of NIH scientists have resulted in major advances. Through the technology transfer process, the promise of "magic bullets" appears closer to realization.

References

1. Frost & Sullivan, Inc. "The U.S. Market for Immunoconjugates Used in Medical Imaging and Therapeutics", Spring 1991, Pages 1-195, Frost & Sullivan, Inc., New York
2. Ehrlich, P., (1913), "Chemotherapy" Proc. 17th Int. Congr. Med. In "The Collected Papers of Paul Ehrlich" (1960) (F. Himmelweit, Ed.), Vol. III. p.510, Pergamon Press, Oxford.
3. Pastan, I., Adhya, S., and FitzGerald, D., U.S. Patent 4,892,827, issued January 9, 1990.
4. Cawley, D., Herschman, H., Gilliland, D. and Collier, R., *Cell*, 22, 563-570 (1980).
5. Lord, J., Roberts, L., Thorpe, P., and Vitetta, E., *Trends in Biotechnology*, 3, No.7, 175-179 (1985).
6. Kohler, G. and Milstein, C., *Nature*, 256, 495-497 (1975).
7. Spooner, R., and Lord, J., *Tibtech*, 8, 189-193 (1990).
8. Blakey, D., Wawrzynczak, E., Wallace, P, and Thorpe, P., in "Monoclonal Antibody Therapy, Prog Allergy, (H. Waldmann, Ed.), Vol 45. pp.50-90, Basel, Karger,(1988).
9. Covell, D., Barbet, J., Holton, O., Black, C., Parker, R., and Weinstein, J., *Cancer Res.*, 46, 3969-3978 (1986).
10. Kaliner, M., and Boltansky, H., U.S. Patent 4,902,495 issued February 20, 1990.
11. Morrison, S., Johnson, M., Herzenberg, L., and Oi, V., *Proc. Natl. Acad. Sci. USA*, 81, 6851-6855 (1984).
12. Reichmann, L., Clark, M., Waldmann, H., and Winter, G., *Nature*, 332, 323-327 (1988).
13. Uhr, J., *Journ. Immunol.* 133, 1 (1984).
14. Bjorn, M., Ring, D., and Frankel, A., *Cancer Res.* 45, 1214-1221 (1985).
15. Greenfield, L., Johnson, V., and Youle, R., *Science* 238, 536 (1987).
16. Hwang, J., FitzGerald, D., Adhya, S., and Pastan, I., *Cell* 48, 129-136 (1987).
17. Jinno, Y., Chaudhary, V., Kondo, T., Adhya, S., FitzGerald, D., and Pastan, I., *Journ. Biol. Chem.* 263, 13203-13207 (1988).
18. Batra, J., Jinno, Y., Chaudhary, V., Kondo, T., Willingham, M., FitzGerald, D., and Pastan, I., *Proc. Natl. Acad. Sci. USA* 86, 8545-8549 (1989).
19. Batra, J., FitzGerald, D., Gately, M., Chaudhary, V., and Pastan, I., *Journ. Biol. Chem.* 265, 15198-15202 (1990).
20. Batra, J., FitzGerald, D., Chaudhary, V., and Pastan, I., *Mol. Cell. Biol.* 11, 2200-2205 (1991).
21. Brinkmann, U., Pai, L., FitzGerald, D., Willingham, M., and Pastan, I., *Proc. Natl. Acad. Sci. USA* 88, 8616-8620 (1991).
22. Rybak, S., Saxena, S., Ackerman, E., and Youle, R., *Journ. Biol. Chem.* 266, 21202-21207 (1991).
23. Rybak, S., Hoogenboom, H., Meade, H., Raus, J, Schwartz, D., and Youle, R., *Proc. Natl. Acad. Sci. USA* 89, 3165-3169 (1992).

AUTOMATED SYSTEM FOR EARLY BREAST CANCER DETECTION IN MAMMOGRAMS

Isaac N. Bankman*, Dong W. Kim*, William A. Christens-Barry*, Irving N. Weinberg#,
Olga B. Gatewood#, and William R. Brody#

The Johns Hopkins University
*Applied Physics Laboratory
#Radiology Department

N 9 3 - 2 5 5 6 8

150 477 -

P. 9

ABSTRACT

The increasing demand on mammographic screening for early breast cancer detection, and the subtlety of early breast cancer signs on mammograms, suggest an automated image processing system that can serve as a diagnostic aid in radiology clinics. We present a fully automated algorithm for detecting clusters of microcalcifications that are the most common signs of early, potentially curable breast cancer. By using the contour map of the mammogram, the algorithm circumvents some of the difficulties encountered with standard image processing methods. The clinical implementation of an automated instrument based on this algorithm is also discussed.

INTRODUCTION

The most prevalent cause of cancer in women is breast cancer. In the United States, one in nine women develops breast cancer in her lifetime and every year more than 170,000 new cases are diagnosed. The incidence of breast cancer is more than double that of colorectal cancer, the second major type in women. However, breast cancer is not the major cause of cancer deaths in women. Studies have indicated that early diagnosis and treatment may significantly improve the 5-year survival for breast cancer patients [1,2]. The American Cancer Society recommends a baseline mammogram for all women by the age of 40, a mammogram approximately every other year between the ages of 40 and 50, and yearly mammogram screening after the age of 50. It has been shown that these screening tests contribute to earlier diagnosis and treatment of breast cancer and many insurance carriers have agreed to cover these examinations. Because awareness and willingness for prevention of breast cancer is increasing rapidly, it is possible that mammography will soon be one of the highest volume X-ray procedures that radiology clinics use regularly. In the U.S. today, about 35 million women are older than 50 and in the next several years, the female U.S. population above the age of 50 will increase at a higher rate than before, reaching about 40 million in the year 2000. While the volume of mammograms is expected to increase, many hospitals are decreasing the number of radiology trainees due to budgetary cuts. The well-recognized goal of performing mammography on a larger scale is becoming more difficult to attain due to the lack of trained readers. Furthermore, the economic feasibility aspects of mammographic screening require that more than 50 mammograms per machine be interpreted daily. This volume is far beyond the current capacity of most mammography clinics in the U.S.

Besides the volume problem, mammographic screening has also an interpretation reliability problem due to the subtlety of the early signs of breast cancer. The life of a women can be saved only if breast cancer can be detected at a very early stage. Early detection of breast cancer in mammograms is a subtle pattern recognition problem due to the wide variation in the normal breast tissue, the large variety of radiographic findings associated with breast cancer, and the similarity between early breast cancer signs and some normal tissue structures. One of the widely used early mammographic indicators of breast malignancy is the presence of clustered microcalcifications. An individual microcalcification appears as a bright spot that ranges in size from about 0.1mm to 2mm in a mammogram. In common mammographic practice, the presence of three or more microcalcifications in a small region (less than 1 cm²) is usually accepted as a cluster. The cluster of microcalcifications is a highly sensitive sign and in many cases it is the first and only sign of an early, potentially curable breast carcinoma [3-5].

With increasing pressure on throughput and the subtlety of early breast cancer signs, the possibility of observer error increases. Fatigue from reading excessive numbers of mammograms contributes to an increase in the number of missed breast cancers [6-8]. Experienced radiologists are aware of the human factors that limit

reliability and they generally stop interpreting mammograms after they have read a certain number on the same day. A reliable, computerized system could contribute much needed speed and accuracy to mammogram interpretation by serving as an assistance device for the radiologist. A computer-aided interpretation tool that indicates suspicious structures in mammograms can allow the radiologist to focus rapidly on the relevant parts of the mammogram. Furthermore, with high resolution film digitization and wide dynamic range, it may be feasible to detect lesions that might otherwise be missed by the radiologist due to their small size or low contrast. An automated system that can recognize reliably early signs of breast cancer and work continuously without fatigue will be a valuable asset for any radiology clinic. Such a system can contribute not only to the availability of vital health care but it has a potential for reducing its cost as well.

In the early breast cancer detection problem, the main goal is to miss as few signs as possible on the mammogram; false negatives can delay diagnosis and preclude the possibility of timely intervention to save the life of the patient. At the same time, false positives are also undesirable because they can cause unwarranted biopsy examinations. Since biopsy requires surgery of the breast, false positives should be minimized. Therefore, in breast cancer detection, both sensitivity and specificity are important, with sensitivity bearing a more vital implication.

Several algorithms have been suggested for detecting clusters of microcalcifications [9-12]. An elegant pioneering algorithm was based on preprocessing of the mammogram for enhancing microcalcifications [9]. This algorithm can be adjusted to provide a sensitivity of 95% or more but introduces 5 or more false positive clusters per mammogram at these sensitivity levels. As stated in the conclusions of [9], to reduce the number of false positives better signal extraction techniques are necessary. A large number of false positives per mammogram that need to be ruled out by the radiologist would cause an undesirable burden in a busy radiology clinic.

In other algorithms such as [10], adequate detection relies on the human operator who has to set manually 10 acceptance thresholds that may vary for different mammograms. Although the image processing and pattern recognition aspects of these algorithms may be effective, they are not directly applicable in a clinical setting due to the human supervision that they require.

Algorithms that use local thresholds derived from the local distribution of intensity values on the mammogram such as [11] rely on the existence of bimodal distributions in local analysis windows where microcalcifications (signal) and normal tissue (noise) form two distinct Gaussian modes. In most cases, the intensity distribution within analysis windows is unimodal and the detection thresholds of this approach are difficult to determine.

A recent approach [12] used clinical information such as age, relatives with breast cancer, biopsy history, breast size, and breast density, combined with shape measurements of microcalcifications using an expert system. This approach yielded 72% accuracy in identifying clusters of microcalcifications.

The high false positive rate of algorithms that use image enhancement may be due to the spectral overlap between signal and noise in the breast cancer detection problem. Because microcalcifications are similar to other small structures in normal tissue as well as small film artifacts, the spatial frequency content of microcalcifications overlaps considerably with that of some normal tissue structures and that of film artifacts. Image enhancement is essentially band-pass filtering in frequency domain and when the spectra of signal and noise overlap to a large extent, the pass band enhances both signal and similar noise components giving rise to false positives. Too little enhancement can preclude the detection of some microcalcifications while too much enhancement can increase significantly the amplitude of small background structures and produce a large number of false detections. The best compromise may change from image to image and can be difficult to determine. Especially when a single enhancement filter is used to enhance all mammograms, poor detection results can be obtained in many cases. This is due to the fact that both microcalcifications and normal tissue structures exhibit a large variability in size and shape in different mammograms. Consequently, the spectra of signal and noise can vary significantly across mammograms. The theoretically optimal approach that has not been applied to mammogram analysis, is to use Wiener filtering [13] that maximizes the correct detection rate. This approach would provide the best band-pass filter for each mammogram, based on prior knowledge of signal and noise spectra. However, two concerns are valid about this approach in the breast cancer detection problem. First, the need for prior knowledge of signal and noise spectra implies a relatively high human guidance for each mammogram where segments of signal and noise have to be indicated to the algorithm. Second, the high level of overlap between signal and noise within the same mammogram undermines the performance of all band-pass filtering approaches including the Wiener filter.

Furthermore, enhancement may introduce an additional difficulty in the development of an appropriate algorithm due to the modification that filtering imparts to the data. The goal of breast cancer detection algorithms is to approximate as closely as possible the recognition performance of experienced radiologists possibly using confirmation by a biopsy examination. Therefore, the target clusters are indicated by radiologists who also provide guidance on the related detection criteria. The interaction with experienced radiologists is essential for the development of a reliable breast cancer detection algorithm. When the image is filtered, in many cases the data used by the algorithm can be considerably different than the data used in visual radiographic interpretation. In such cases the detection criteria and suggestions of the radiologist may not be directly applicable to the algorithm, and consequently the accordance between visual and automated detection decreases.

From the information theory point of view, if the mammogram is digitized appropriately, all the information needed to detect microcalcifications is present in the raw image. Enhancement is an attempt to eliminate irrelevant and obscuring information and to transform the relevant information for more convenient detection. Since all the information needed is in the raw data, it is possible that algorithms that can access the relevant information without enhancement can be developed.

In most available algorithms for breast cancer recognition, the detection is performed by comparing the amplitude of the signal, i.e. the local intensity of the mammogram to a threshold. In difficult pattern recognition problems where the signal and noise are similar in spectral content as well as in amplitude, successful detection has been achieved by extracting relevant features from the data [14] while detection algorithms based only on amplitude performed poorly [15]. Especially when the goal is to approximate the visual interpretation of the data, features that reflect the visual cues convey the most relevant and effective information. Similarly, in mammogram analysis the visual recognition criteria developed by expert radiologists across many years can guide the development of an effective algorithm by suggesting features that characterize microcalcifications. An additional advantage of features that represent visual cues is that they provide a set of parameters that can be easily interpreted. This allows a more effective interaction with radiologists and gives the algorithm a potentially higher degree of acceptance in the radiology community.

In some of the available algorithms for breast cancer recognition, estimates of the local intensity gradient are used for detection because microcalcifications have a relatively higher intensity with respect to their immediate surroundings. This is done by comparing the pixel values within a small square kernel about the size of a microcalcification, to the pixel values outside the kernel. Because square kernels do not match the shapes of microcalcifications adequately, these estimates of local gradients can be misleading. In fact any measurement for characterizing microcalcifications may be inadequate if it is made by observing the interior of a kernel of preset arbitrary shape and size.

Based on the considerations mentioned above, we set the following specifications for the design of a new algorithm:

1. Operation on raw data without enhancement.
2. Use of features representing visual mammogram interpretation criteria.
3. Operation without preset analysis kernels.
4. Operation without assumptions about the statistical distribution of parameters.
5. Completely automated operation without human intervention.

The algorithm that we developed satisfies these specifications and circumvents some of the difficulties encountered in other algorithms.

DATA

The data were obtained by digitizing 9 mammograms from different patients diagnosed to have cancer by radiographic examination as well as biopsy. Each mammogram was annotated by an experienced radiologist who indicated the locations of all clusters of microcalcifications in the mammograms. A total of 13 clusters were annotated.

Mammograms were backlit using a uniform source light box and digitized in overlapping segments of 25.6 mm height by 38.4 mm width. Segments were overlapped by about 20% in each dimension, eliminating the

possibility that a microcalcification might appear on a segment boundary and escape detection during numerical analysis. Each segment was imaged by a Canon FD 50 macro lens (with extension tube) onto a Sony XC-77ce CCD array camera at a spatial resolution of 50 μm . The illumination intensity was adjusted so that saturation did not occur in any of the signal bearing regions of the mammograms. The data from each segment, consisting of a raster array of 512 by 768 pixels of 8-bit gray scale, were stored on magnetic media for subsequent numerical analysis.

DEVELOPED ALGORITHM

The strategy of the detection algorithm is to view the image as a landscape where elevation corresponds to brightness. In this perspective, microcalcifications appear as prominent peaks that stand out with respect to the local surround. A section of a mammogram that contains a microcalcification cluster is shown in Fig. 1 and the corresponding 3-D plot of the cluster is shown in Fig. 2a. The algorithm starts by forming the contour map of the image. The contour plot obtained in the vicinity of the cluster is shown in Fig. 2b. These contours are iso-intensity contours analogous to iso-elevation contours in cartography and therefore, they are *not* obtained by edge detection and *do not* require local gradient estimates.

When the contours are obtained, the detection algorithm focuses on concentric contours. Each set of concentric contours that represents a peak (an individual microstructure) is analyzed separately. From each peak, the algorithm obtains a sequence of contour areas progressing from the highest contour level in that set (small area) towards low contour levels (larger area). Contour areas that are too small or too big to be part of a microcalcification are not included in the area sequence of a peak. The algorithm is designed to determine the area growth sequence of an individual peak when other peaks are close, by accounting for the merging of contours.

The algorithm computes 5 measurements (features) from the area sequence of peaks:

- 1) **Departure.** In visual inspection, microcalcifications are bright structures with a relatively sharp appearance in their visually perceived edge. In the landscape view of a digitized image, the perceived sharpness of a microstructure depends on the departure of that peak from the surrounding background. A microstructure with sharp edges is a peak that departs abruptly from the background while a fuzzy microstructure is a peak that departs very gradually from the surrounding background. The departure feature quantifies the sharpness of a microstructure using the area sequence of that peak. In the area sequence of a peak, an abrupt departure from background is reflected as a sudden change in the rate of change of the area sequence near the base of the peak. This information can be obtained by using the second derivative of the area sequence. In order to obtain a departure value that is insensitive to the size of the microstructure (absolute values of the areas), the algorithm computes the first derivative sequence, and sets the departure to the maximal relative change in the first derivative, in the lower half of the peak. The contour level where the departure is obtained is considered the base of the peak, i.e. the immediate background level.
- 2) **Prominence.** This parameter reflects the relative brightness cue that is used in visual inspection. This local contrast information contributes to discriminating microcalcifications from both normal tissue structures and film artifacts. The prominence value is set to the number of contours above the departure level and it is approximately proportional to the brightness difference between the brightest region of the microstructure and the immediate surround at the level of departure from the background.
- 3) **Steepness.** In addition to the sharpness at the perceived edge which is reflected by the departure feature, the rate of change of intensity throughout a microstructure is a significant property for visual inspection. Generally, normal breast tissue structures appear globally more diffuse than microcalcifications. Such diffuse structures are represented by peaks that have a gradual increase in height. In the landscape view of the mammogram, peaks that correspond to microcalcifications have a higher overall steepness than normal tissue structures. Moreover, the peaks of some film artifacts are typically steeper than microcalcification peaks. The steepness parameter is obtained by using the first derivative of the area sequence in a manner that results in higher values for steeper peaks.
- 4) **Distinctness.** In many cases, the normal breast tissue in a mammogram has a grainy appearance due to a large number of contiguous normal microstructures. Although microcalcifications may be clustered occasionally in close proximity to each other, they are more distinct and separate from each other as well as from normal

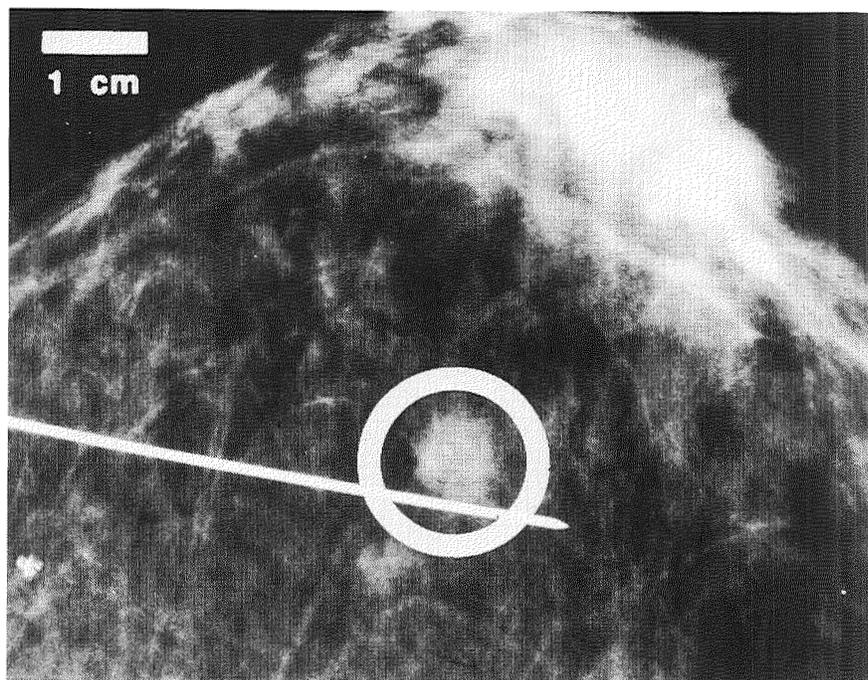
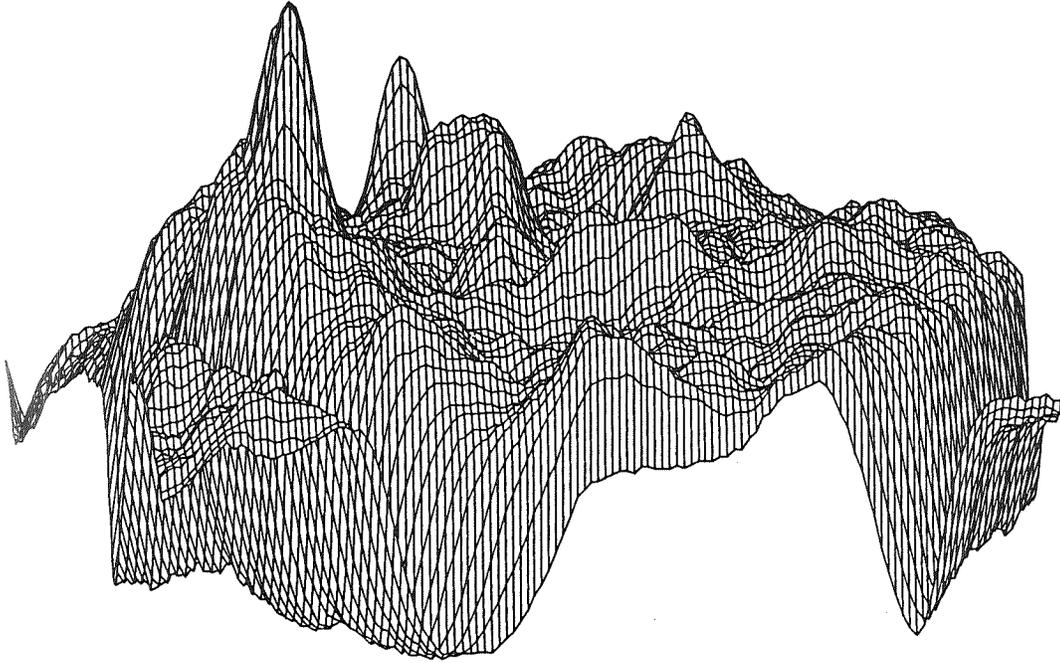


Figure 1. Photographic enlargement of a mammogram analyzed in this study. A microcalcification cluster (circled) is shown.

a



b

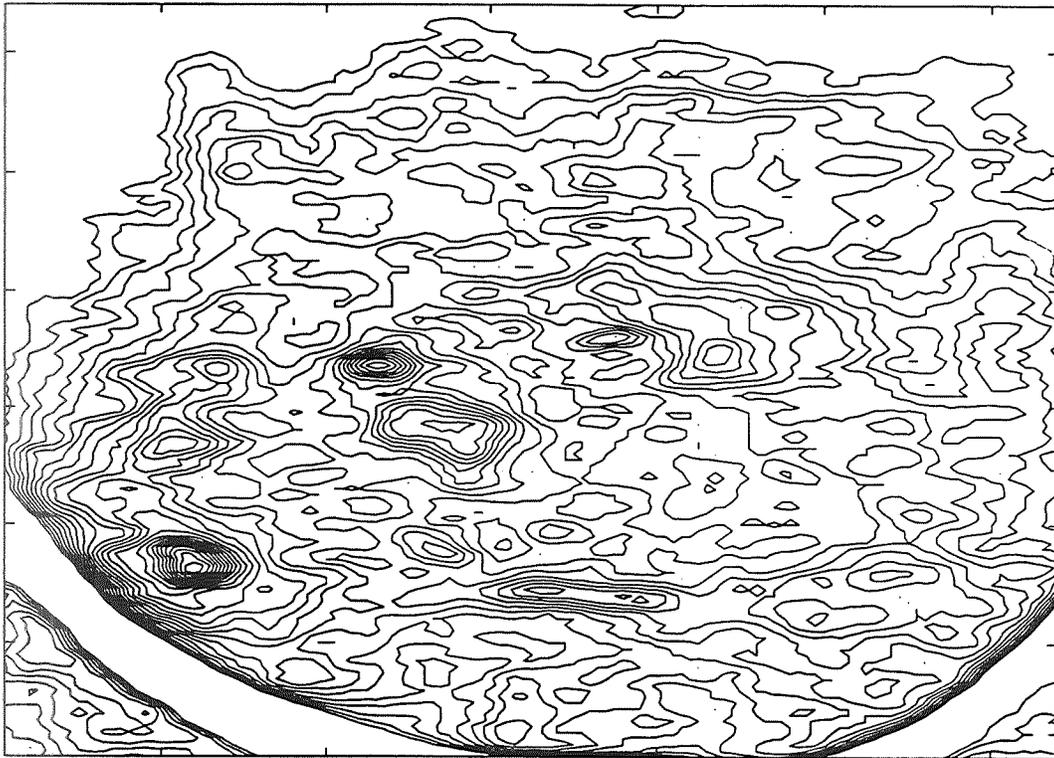


Figure 2. (a) Intensity surface plot of the region containing a microcalcification cluster shown in Fig. 1; (b) Iso-intensity contour plot derived from the same region.

microstructures. The distinctness of a peak is set to the number of contour levels between the tip and the level where its contour merges with that of the nearest peak.

5) Compactness. The edge morphology of a microcalcification is a significant visual cue. For each peak, this morphological information is obtained by using the characteristic contour of a peak obtained just above the merging level. The compactness feature is computed using the ratio of the perimeter to the area of the characteristic contour. Compactness is a standard morphological descriptor that has a value of 1 for a circle and increases as the shape becomes more irregular. The compactness of some types of artifacts and most normal tissue structures is relatively higher than that of microcalcifications.

Each peak is characterized with the 5 features that the algorithm extracts from the raw mammogram data and the discriminant between microcalcifications and other structures is based on these features. The discrimination can be performed with conventional Bayesian classification, standard feedforward neural networks [16], or specialized neural networks [e.g. 17]. In this study the Bayesian classifier was used and adequate results were obtained.

The decision parameters of the classifier were determined on 3 mammograms that formed the training set. The digitized training mammograms formed a data-base of 64 image segments containing more than 1000 microstructures that had a size of interest (less than 2 mm wide). The training set contained 5 of the microcalcification clusters indicated by the radiologist. The distributions of the features were obtained for the populations of microstructures within the indicated clusters (detection class) and for the population of peaks in the rest of the mammogram (rejection class) separately. The decision thresholds were set in order to maximize the discrimination between the detection class and the rejection class. A cluster was indicated by the algorithm when 3 or more microcalcifications occurred in an area of less than 1cm^2 using the 5 features and a two-phase data reduction approach.

EVALUATION

The performance of the algorithm was evaluated on the 6 other mammograms that formed the test set. The digitized test set resulted in 84 image segments containing more than 1200 candidate microstructures. The test set contained 8 of the clusters indicated by the radiologist and the algorithm detected all 8. In addition, the algorithm detected 1 false positive cluster in one mammogram. Therefore, on the test set the sensitivity was 100% with 0.17 false clusters per mammogram.

In pattern recognition applications, the performance of an approach is measured by the balance of false negatives and false positives that it can provide. Almost any algorithm can be made sensitive enough to detect all events of interest (no false negatives). However, increasing the sensitivity generally reduces the specificity and causes a larger number of false positives. Therefore, in many pattern recognition applications the false positive rate associated with a desirable sensitivity level is used as a measure of performance. A sensitivity level of about 95% or more is desirable in early breast cancer detection. For such a high sensitivity level, the 0.17 false clusters per mammogram obtained with this algorithm provide a considerably better specificity than 5 or more false positive clusters per mammogram obtained with other algorithms.

This algorithm was developed specifically to detect microcalcifications based on the radiographic visual evaluation criteria. These criteria were computationally expressed as features extracted from the raw data without using enhancement. The use of the contour plot provided a convenient technique for computing the features without using preset arbitrary analysis kernels. In this manner, all measurements were obtained using natural morphological contours of microcalcification peaks. The decision thresholds were applied to features and not to the intensity data. Appropriate values of these thresholds were determined using a large number of diverse microstructures. Therefore these thresholds did not depend on local statistics, they held across mammograms and did not have to be adjusted for each mammogram separately. Once the thresholds were set using a representative training set, the algorithm operated in a fully automated manner without human supervision. The algorithm will be further validated on more than hundred mammograms during use in the Department of Radiology of The Johns Hopkins Hospital.

CLINICAL IMPLEMENTATION

The automated system based on this algorithm will be a reliable diagnostic tool that can assist radiologists in early breast cancer detection on mammograms. The system will be made of a scanner, a low-cost workstation, a high-resolution display and a printer for hard copies of results. The speed of the workstation will allow one mammogram to be analyzed in less than 5 minutes. The fully automated operation of the algorithm ensures that the system will not introduce an additional burden to radiologists.

The automated system will be able to analyze a mammogram without human supervision; however, it will be designed to benefit from the experience of radiologists across time. This will be possible with training software that will be available to radiologists or technicians. Occasionally, a human operator will indicate to the system the location of false negative or false positive microcalcifications, using the keyboard or mouse. The training software will automatically adjust the operational parameters of the system to detect the microcalcifications that were missed and to reject the false positive structures. This will be achieved with minimal change on past correct performance because the adjustment will be made by taking into account not only the currently indicated structures but an archive of some previously encountered structures. This archive will contain the features of a large number of microcalcifications as well as other structures (normal tissue, artifacts, etc.) that were located very close to the decision boundary between these two classes. Therefore, the structures in the archive will be those that would be affected first by changes in operational parameters. The training software will automatically optimize the discrimination based on the past examples and the currently indicated structures. In this adjustment, the weight given to current structures will be user-selectable.

The automated operation of this system is especially suited for clinics that have to screen a large number of mammograms every day. In such a clinical setting, the system can operate virtually all the time, as long as a technician is available for feeding the mammograms to the scanner. An automated feeding instrument can also be conceived. Assuming 10 hours of operation per day and a worst case of 5 minutes per mammogram, about 120 mammograms a day can be screened by the system without requiring any time from the radiologist. When the results of the automated system are available, the expert radiologist will focus on the regions where the system indicated microcalcification clusters in each mammogram to confirm the results. For the purpose of quality control, the radiologist might also screen some regions that were cleared by the system on several mammograms. The expected clinical benefits are: i) accurate detection of subtle signs of breast cancer that might be missed by radiologists and, ii) significant reduction in the amount of time that radiologists spend for screening mammograms. Currently, due to the subtlety of early breast cancer signs, radiologists use a magnifying glass to screen mammograms. The time required for the visual interpretation of a complete mammogram can often reach 15 minutes and in some cases it can take up to 30 minutes. The automated system is expected to reduce the time required of the radiologist by an order of magnitude.

REFERENCES

1. American Cancer Society, *Ca Cancer J. Clin.* 33, 226 (1982).
2. D.B. Kopans, J.E. Mayer and N. Sadowsky, *N. Engl. J. Med.* 310, 960 (1984).
3. L. Bassett, "Mammographic analysis of calcifications", *Radiologic Clin. N. Amer.*, 30(1), 95 (1992).
4. P.C. Stomper, J.C. Connolly, J.E. Meyer et al, "Clinically occult ductal carcinoma in situ detected with mammography: Analysis of 100 cases with radiologic-pathologic correlation", *Radiology*, 172, 235 (1989).
5. E.A. Sickles, "Mammographic Features of 300 consecutive nonpalpable breast cancers", *Am. J. Roentgenol.*, 146, 661-663 (1986).
6. M.J. Smith, "Error and variation in diagnostic radiology", pub. Charles C. Thomas, Springfield, IL (1967).
7. L. Kalisher, "Factors influencing false negative rates in xeromammography", *Radiology* 133, 297-301 (1979).
8. J.E. Martin, M. Moskowitz, and J.R. Milbrath, "Breast cancer missed by mammography", *Amer. J. Roentgenol.* 132, 913-918 (1982).
9. H.P. Chan, K. Doi, C.J. Vyborny, K.L. Lam and R.A. Schmidt, "Computer-aided detection of microcalcifications in mammograms", *Investigative Radiology* 9, 664 (1988).
10. B.W. Fam, S.L. Olson, P.F. Winter, F.J. Scholz, "Algorithm for the detection of fine clustered microcalcifications on mammograms", *Radiology* 169, 333-337 (1988).
11. D.H. Davies and D.R. Dance, "Automated computer detection of clustered microcalcifications in digital mammograms", *Physics in Medicine and Biology* 35(8), 1111-1118 (1990).

12. E.A. Patrick, M. Moskowitz, V.T. Manshukhani and E.I. Gruenstein, "Expert learning system network for diagnosis of breast calcifications", *Investigative Radiology* 26, 534-539 (1991).
13. R.C. Gonzales and P. Wintz, "Digital Image Processing", 229-232, pub. Addison-Wesley, Reading, MA (1987).
14. I.N. Bankman, V.G. Sigillito, R.A. Wise, and P.L. Smith, "Feature-based detection of the K-complex wave in the human electroencephalogram using neural networks", *IEEE Trans. on Biomed. Eng.* (1992), in press.
15. B.H. Jansen, "Artificial neural nets for K-complex detection", *IEEE Eng. Med. Biol.*, 9, 50-52, (1990).
16. D. E. Rumelhart, J. L. McClelland, and the PDP Research Group, "Learning internal representations by error propagation", MIT press, Cambridge, MA, 318-328, (1988).
17. I.N. Bankman, J. Sadowsky, V.G. Sigillito, "A Neural Network for Adaptive Surfaces and Closed Decision Boundaries", *Proc. 25th Annual Conf. on Information Sciences and Systems*, 25, 848-853, (1991).

58-52
153478
P-10

**THE DESIGN OF MECHANICALLY COMPATIBLE FASTENERS
FOR HUMAN MANDIBLE RECONSTRUCTION**

**Jack C. Roberts, Senior Professional Staff
John A. Ecker, Senior Professional Staff
Paul J. Biermann, Senior Professional Staff**

**The Johns Hopkins University
Applied Physics Laboratory
Laurel, Maryland**

ABSTRACT

Mechanically compatible fasteners for use with thin or weakened bone sections in the human mandible are being developed to help reduce large strain discontinuities across the bone/implant interface. Materials being considered for these fasteners are a Polyetheretherketone (PEEK) resin with continuous quartz or carbon fiber for the screw. The screws were designed to have a shear strength equivalent to that of compact/trabecular bone and to be used with a conventional nut, nut plate, or an expandable shank/blind nut made of a ceramic filled polymer. Physical and finite element models of the mandible were developed in order to help select the best material fastener design. The models replicate the softer inner core of trabecular bone and the hard outer shell of compact bone. The inner core of the physical model consisted of an expanding foam and the hard outer shell consisted of ceramic particles in an epoxy matrix. This model has some of the cutting and drilling attributes of bone and may be appropriate as an educational tool for surgeons and medical students. The finite element model was exercised to establish boundary conditions consistent with the stress profiles associated with mandible bite forces and muscle loads. Work is continuing to compare stress/strain profiles of a reconstructed mandible with the results from the finite element model. When optimized, these design and fastening techniques may be applicable, not only to other skeletal structures, but to any composite structure.

INTRODUCTION

During mandible reconstructive surgery many problems may be encountered when attaching thin or weakened (diseased) sections of bone to one another, or to a replacement material. Some of these problems may stem from differences in the mechanical properties of bone and the implant material, or from the fastening method. The use of bonding techniques may result in premature failure of the interface due to the reduced or weakened bone section at the interface. The use of mechanical fastening techniques could result in shearout of the thin/weak bone due to the higher stiffness of conventional fastener materials, as well as large strain discontinuities across the interfacial boundary between different materials.

Atrophy, due to stress shielding of underlying bone in metal fixation devices, is said to be the most important reason for the removal of rigid metallic plates and screws. Several investigators have proposed the use of biodegradable materials with the same properties as bone to be used in place of metal fixation devices (1-3). These biodegradable materials could be used in situations where the fixation device is no longer needed after the fracture has healed. To create a better interfacial bond between implants and bone other investigators have explored the use of porous metallic implant materials (4), and to eliminate corrosion porous polysulfone (5) and porous hydroxyapatite (6,7). Hydroxyapatite was thought to create a better bond with bone because it is one of the constituents of bone. However, in situations where the implant or reconstructed mandible must remain in place with the fixation device, it is desirable to have a fixation device that has the same properties as bone and will have adequate strength at the bone/implant interface.

Therefore, the purpose of this study was to design non-metallic fasteners having properties similar to bone for a reconstructed mandible or mandibular implant. The fastening method should "design in" the ability of the fastener to flex with bone, thus preventing bone atrophy, while providing a continuous load path from bone to the replacement. The design of the fasteners will be guided by both a physical model and a finite element model of the mandible.

MATERIALS AND FASTENER DESIGN

Materials

Bone has a hard outer surface (compact bone) and a soft, porous inner core (trabecular bone). This makes it similar to sandwiched composite structures used in the aerospace field. These structures consist of an highly porous inner core (usually a honeycomb material to carry shear loads) sandwiched between aluminum, glass/epoxy or graphite/epoxy face sheets (to carry the tensile or compressive loads). A comparison of the two assembled structures is shown in Figure-1. Because of the similarity in function, the implant designs will use fasteners similar to those used with sandwich composites in the aerospace field. In order to size the screws in any of the following designs, bone properties in shear (8) were used to calculate the necessary screw diameter. In order to minimize strain discontinuity between bone and fastener, the fastener will be designed to have elastic properties equivalent to that of bone. Since the fasteners will have a smaller cross sectional area than that of bone, the shear strength will need to be higher than that of bone. The screw materials that most closely met these criteria were Polyetheretherketone (PEEK) with either continuous carbon or continuous quartz fibers¹. This design methodology may be a better way to select fasteners for any composite structure, not just bone.

Fastener Design

For ease of assembly, the fastener design selected should allow the surgeon to insert the fastener from the buccal (outside) surface of the mandible. To reduce trauma to the patient and prevent possible infection, the fastener should be flush with the surface of the replacement mandible or mandibular bone. It may also be advantageous to have the fastener put the bone into compression to prevent atrophy. All of these design "criteria" need to be considered when selecting the optimum fastener design.

Three fastener designs are shown in Figure-2. Each design relies on a screw and a nut or nut plate to retain the mandibular replacement. The design in Figure-2A, shows a screw assembled from the lingual (inside) of the mandible into a blind hole in the replacement mandible. The underside of the head of the screw is contoured so that the assembly load is gradually spread over compact bone on the mandible's buccal (outside) surface. The prosthesis or implant has internally threaded blind holes. This design requires the use of a screw driver from the lingual side of the mandible. The design in Figure-2B, uses blind nuts installed from the lingual side of the mandible. A screw is inserted through holes in the buccal side of the replacement mandible. This design requires the blind nuts to be countersunk from the lingual side. The third design (Figure-2C) again features a screw countersunk into the replacement mandible from the buccal side. This design is unique, in that, the nuts have left hand threads on the outside and right hand threads on the inside. Thus, during assembly, no counterboring is required on the lingual side. The nuts will self lock from the buccal side by the use of a simple tool. In all of these designs a jig could be used to drill holes into mandibular bone that correspond with those predrilled in the replacement mandibular or nuts. The shear and tensile strength of the screws in any of the aforementioned designs should exceed those of bone. The compressive modulus should match that of bone to prevent strain discontinuities between the replacement and bone.

An alternative design that relies on an interference fit between a screw and a self-clinching expansion nut rather than tension as in the previous designs, is the shown in Figure-3. The self-clinching expansion nut would be made of ceramic particulate filled polymer. The nut has scores or flutes on the outer surface to allow it to separate under load. An oversized screw is inserted into an undersized tapped hole in the nut. The wall is sized such that the interference fit of the screw and cylinder causes the cylinder walls to expand and split. The wall segments are pushed outward and trapped between the fastener and compact/trabecular bone. The polymer used in this cylinder could be a thermoset or thermoplastic resin filled with enough ceramic to create a somewhat brittle material, but not brittle enough to crush under the screw compressive load. If a thermoplastic is used, enough filler would have to be added to prevent creep under compressive load. The outside surface of the cylinder could be tapered in such a way as to optimize the pressure profile on the outer surface of the fastener. The neck of the fastener could be sized to allow the driver head to be torqued off when the proper preload torque is reached. The entire installation requires only boring holes and no fastener is required on the lingual side of the mandible. The fastener design from Figure-2B or 2C is shown in Figure-4 as an example of a partial replacement mandible.

¹ CPN800A-06-03 (PEEK/Long Carbon), CPN800J-06-03 (PEEK/Long Carbon), Cherry Textron, Santa Ana, CA.

PHYSICAL AND FINITE ELEMENT MODELS

Physical Model

A physical model of the mandible has been developed and fabricated to help with selection of the optimal fastening technique. The fabrication process involves a two step molding process. The first mold is used to fabricate the replacement for trabecular bone. Once solidified, this structure is positioned inside the second mold where a substitute for compact bone is formed over it. Once completed the external geometry of the model, formed by the mold, duplicates that of a human mandible. After evaluating many combinations, the materials that gave cutting and drilling properties similar to bone, were a foamable polymer for trabecular bone, and a ceramic filled epoxy for compact bone. The replacement mandible is shown in Figure-5 and a section through the replacement mandible is compared to a section through a human mandible in Figure-6. This physical model was drilled and cut using medical drills and saws, i.e., a Synthes drill² and a standard Micro-E sagittal saw. The replacement mandible in a dry state, cut and burned like bone. However, while under irrigation the saw blade bound in the material like it would in human bone. The replacement mandible drilled and tapped similarly to bone.

In addition to its use as a tool in the selection of the best fastening technique, the replacement mandible could also, and perhaps more importantly, be used as a training tool for surgeons and medical students.

Finite Element Model

PDA-PATRAN³ was used to create a finite element model of 1/2 the mandible as shown in Figure-7. This model consists of 7560 nodal points and 6716 solid isoparametric hexagonal elements with symmetry boundary conditions applied at the mid-plane and fixed in the condylar region. A very refined mesh was used along the buccal/lingual side in order to provide for easy modifications for incorporation of fastener devices. Both trabecular and compact bone properties were used in the model as taken from Ref. (9). Figure-8 shows stress contours in the mandible model for a static analysis incorporating a first molar point load of 250 N (1,100 lb) and appropriate muscle forces as taken from Ref. (10). MSC/NASTRAN⁴ and COSMOS/M⁵ will be used to exercise this model with different loading, boundary conditions and alternate material properties. The results will be compared to the physical model under similar loads. Once the finite element model has been verified, it can be used to determine the optimal fastening technique.

SUMMARY AND CONCLUSIONS

Several fastener designs have been proposed as fastening devices for a replacement mandible. The first design simply relies on a composite screw inserted into the human mandible on the lingual surface and threaded into a blind hole in the replacement mandible. The second design consists of a composite screw countersunk into a replacement mandible on the buccal surface and threaded into a nut countersunk into the human mandible. The third design relies on right hand threads on the inside and left hand threads on the outside of a nut on the lingual surface of a human mandible to secure a composite screw whose head is countersunk and locked into the replacement mandible. The last design relies on the interference fit between an oversized composite screw and an undersized hole in a ceramic filled epoxy self-clinching expansion nut to force the ceramic to expand outward into the human mandible. These designs must be analytically modeled and experimentally verified before a final device is selected.

² Synthes Ltd. USA, Wayne, PA.

³ PDA-PATRAN finite element pre-and post-processor code, PDA Engineering, Costa Mesa ,CA.

⁴ MSC/NASTRAN finite element solver, MacNeal-Schwendler Corp., Los Angles, CA.

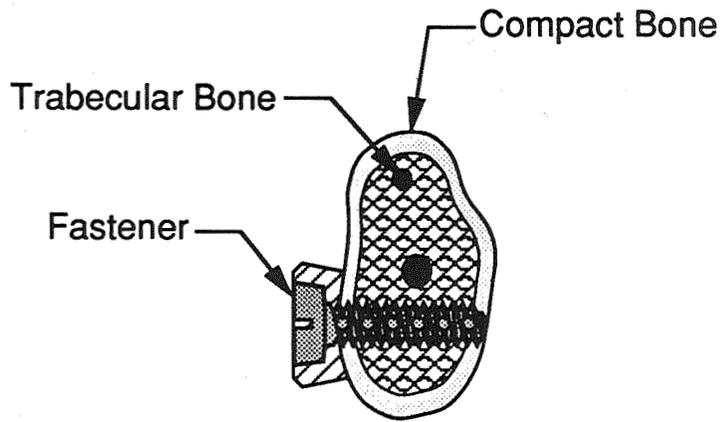
⁵ COSMOS/M finite element software, Structural Research and Analysis Corp., Santa Monica, CA.

A physical model of the mandible was fabricated that included both trabecular and compact bone substitutes. An expanding foam was used to represent trabecular bone and a ceramic filled epoxy was used to simulate compact bone. When cut and drilled, this material acted similarly to bone. Besides being used to test the fastener designs, this substitute mandible could be of benefit to surgeons and medical students training in orthopedic surgery.

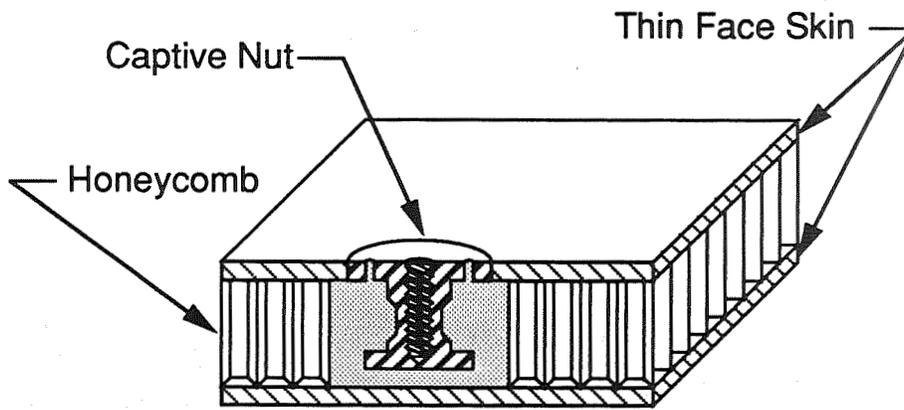
The mandible has been modeled using the finite element technique. Several different loading conditions will be applied to both the physical and finite element models and their results compared. After the finite element model is verified, it will be used to evaluate stress states in each of the fastening techniques. Once optimized for the specific design criteria, the selected fastening technique should not be limited to the mandible or other skeletal structures, but also to composite structures in general.

REFERENCES

1. Lewis, D.H., Dunn, R.L. and Casper, R.A., "Development Of Biodegradable Implants For Use In Maxillofacial Surgery", Southern Research Institute, SORI-EAS-82-507, June 1982.
2. Bos, R.R.M, Rozema, F.R., Boering, G., Nijenhuis, A.J., Penning, A.J. and Jansen, H.W.B., "Bone-Plates And Screws Of Bioabsorbable Poly (L-Lactide) - An Animal Pilot Study", British Journal of Oral and Maxillofacial Surgery, 27(6), 1898, 467-476.
3. Ibay, A.C., "Development Of A Moldable, Resorbable Appliance For Use In Maxillofacial Surgery", Southern Research Institute, SRI-APC-91-571, June 1991.
4. Cook, S.D., Klawitter, J.J. and Weinstein, A.M., "Model For The Implant-Bone Interface Characteristics Of Porous Dental Implants", J. Dental Research, August, 1982, 1006-1009.
5. Ballintyn, N.J. and Spector, M., "Porous Polysulfone As An Attachment Vehicle For Orthopedic And Dental Implants", Biomat., Med. Dev., Art. Org., 7(1), 1979, 23-29.
6. Schliephake, H. and Neukam, F.W., "Bone Replacement With Porous Hydroxyapatite Blocks and Titanium Screw Implants: An Experimental Study", J. Oral Maxillofacial Surgery, 49(2), 1991, 151-156.
7. Denissen, H.W., Kalk, W., de Nieuport, H.M., Maltha, J.C. and van de Hooff, A., "Mandibular Bone Response To Plasma-Sprayed Coatings Of Hydroxyapatite", Int. J. Prosth., 3(1), 1990, 53-58.
8. Reilly, D.T. and Burstein, A.H., "The Mechanical Properties of Cortical Bone", J. Bone and Joint Surgery, 56A(5), July, 1974, 1001-1022.
9. Hart, R.T., Hennebed, V.V., Thonogreda N., VanBuskirk, W.C. and Anderson, R.C., "Modeling The Biomechanics of The Mandible: A Three-Dimensional Finite Element Study", J Biomechanics, 25(3), 1992, 261-286.
10. Haskell, B., Day, M. and Tetz, J., "Computer-Aided Modeling In The Assessment Of The Biomechanical Determinants Of Diverse Skeletal Patterns", Am. J. Orthodontics, 89(5), 1986, 363-382.



Bone Architecture and Current Fastening Technique



Aerospace Fastener for Honeycomb Panels

FIGURE-1 Comparison of Fastening in bone to That in Aerospace Composites

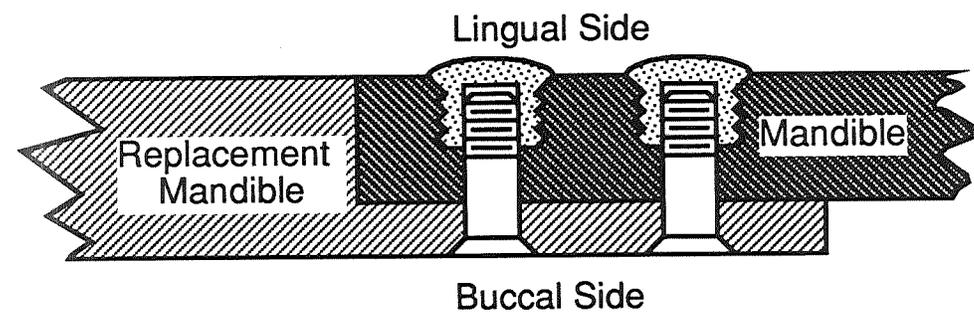
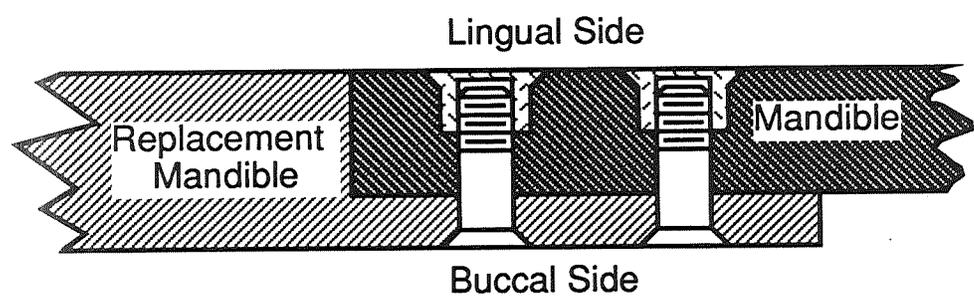
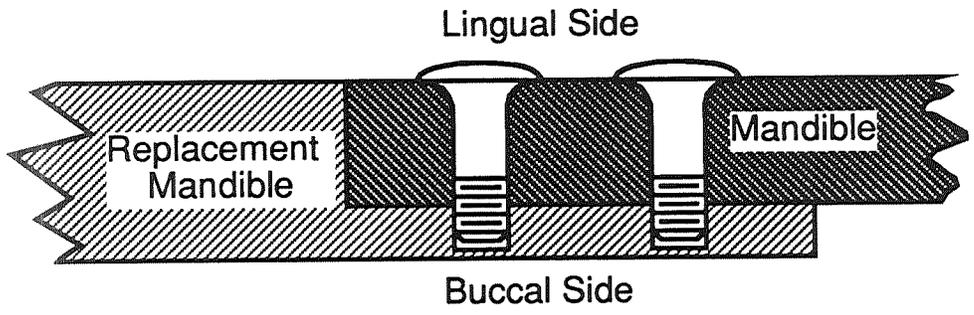


FIGURE-2 Proposed Fastener Designs for Replacement Mandible

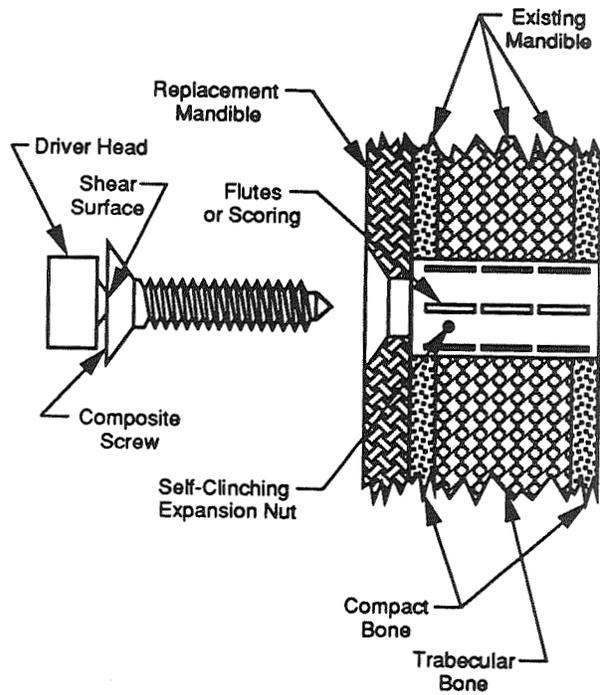


FIGURE-3 Alternative Fastener Design That Relies on an Interference Fit

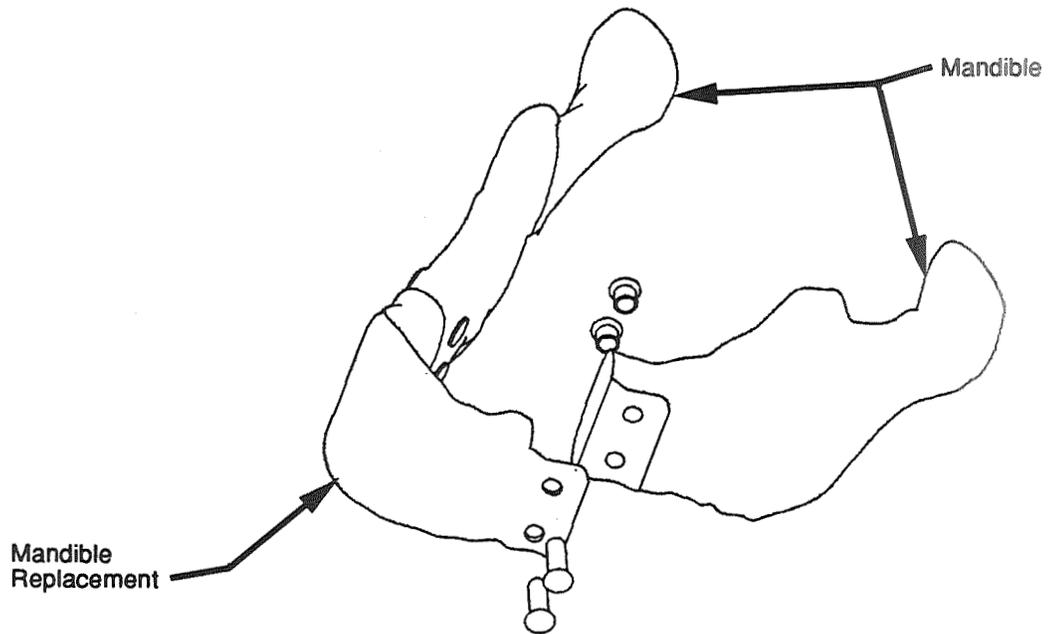


FIGURE - 4 Partial Replacement Mandible With PEEK Composite Fasteners



Figure - 5 Replacement Mandible

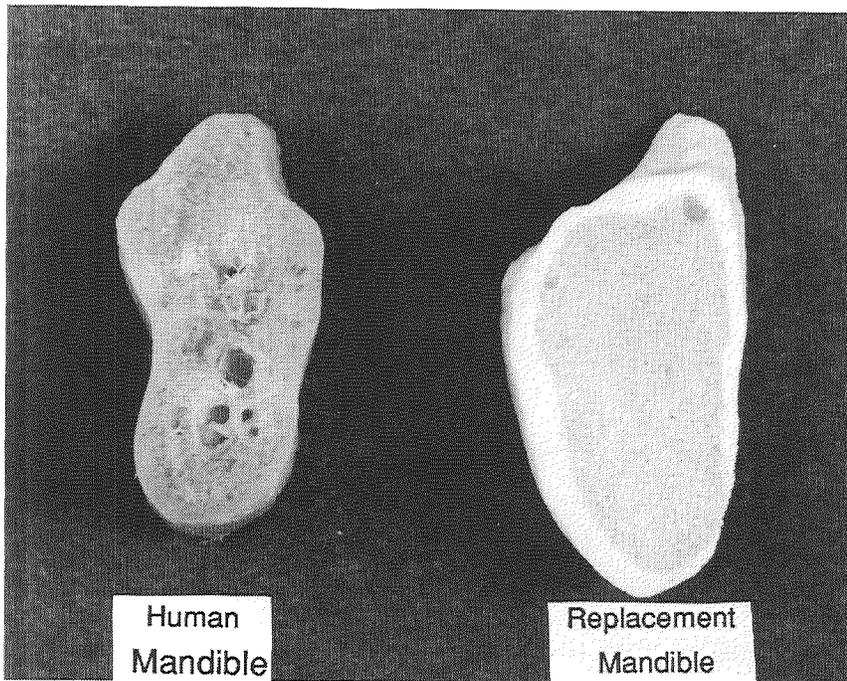


Figure - 6 Section Through Replacement Mandible Compared to Section Through Human Mandible

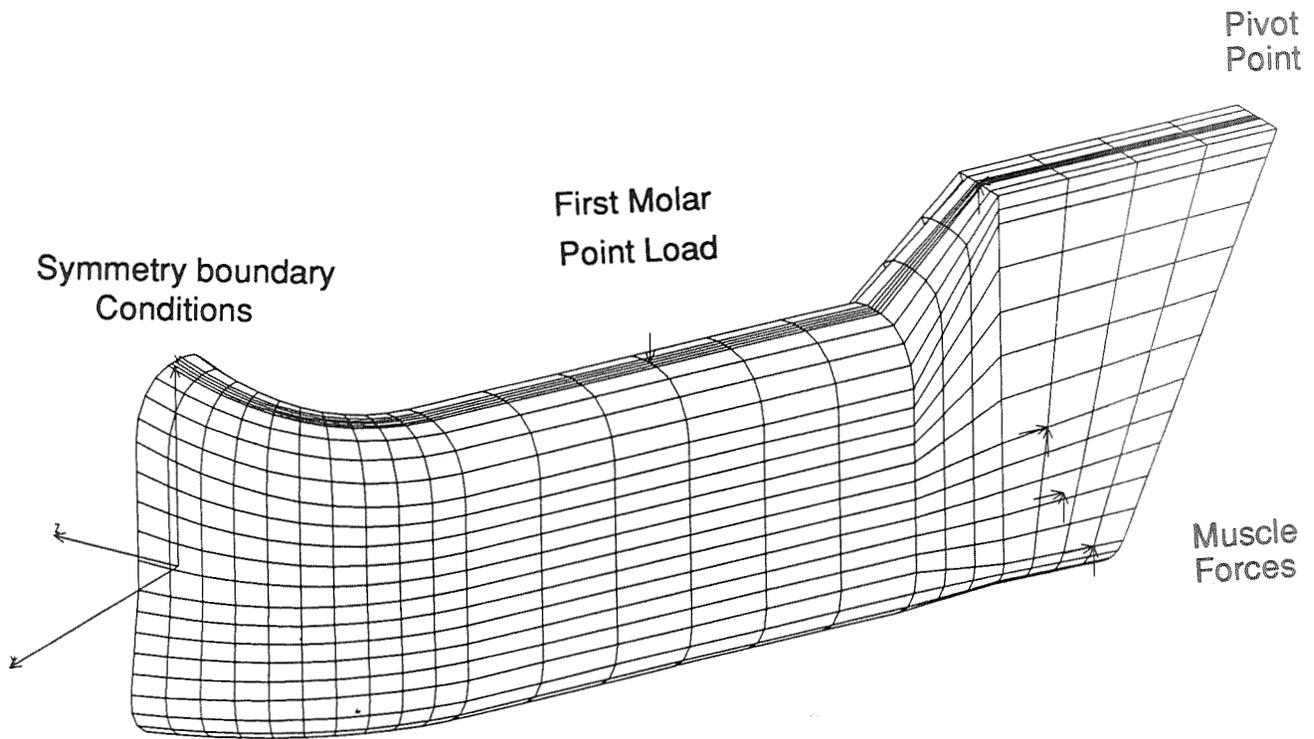


Figure - 7 Finite Element Model of 1/2 The Human Mandible

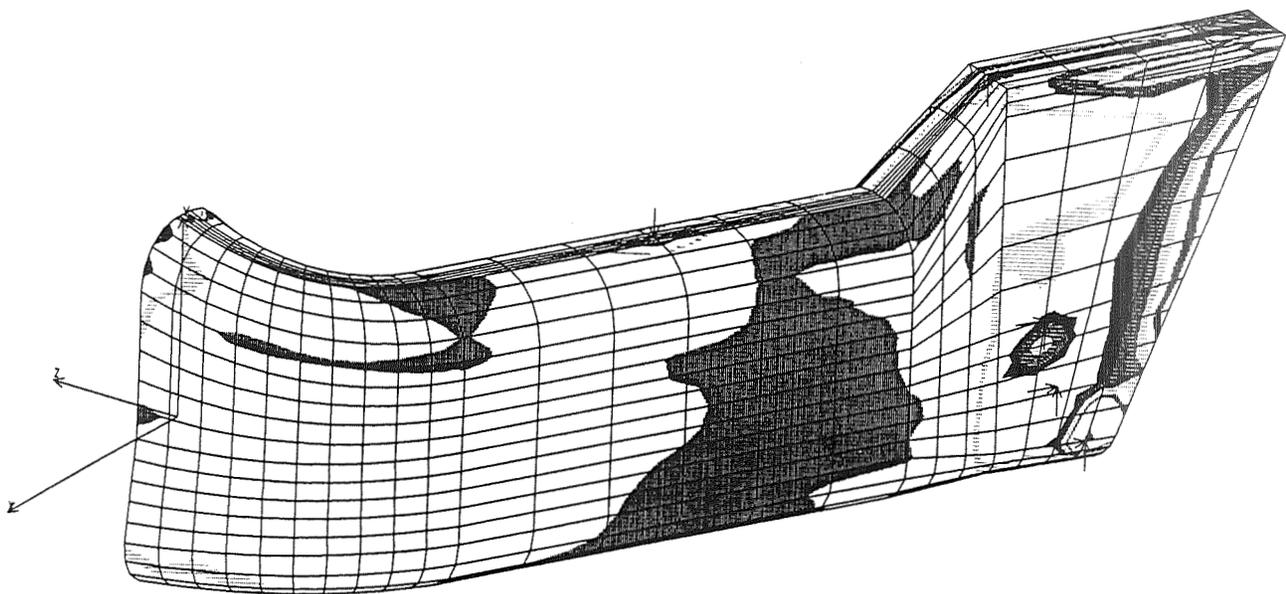
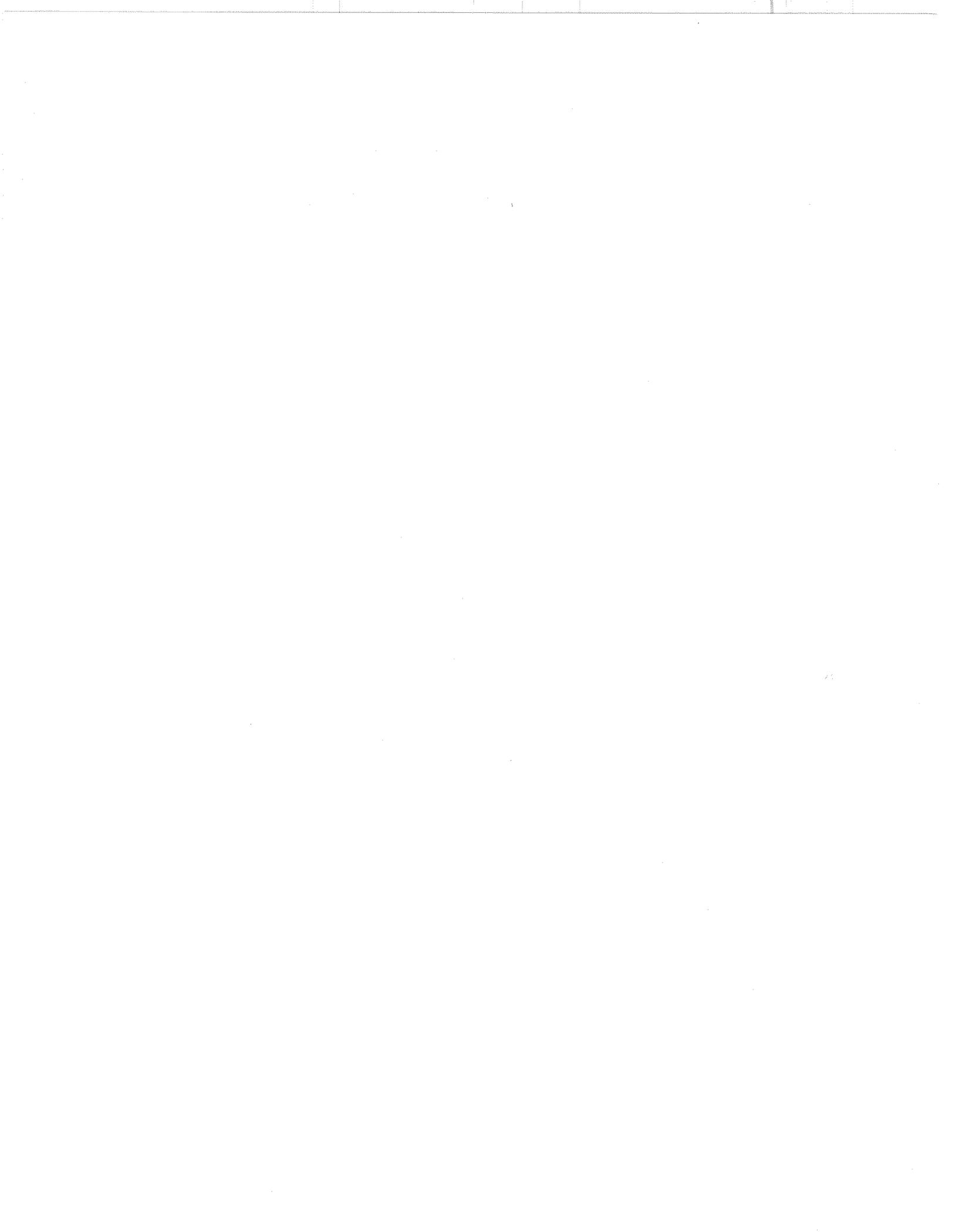


Figure - 8 Stress Contours in a Human Mandible Under a 250 N (1,100 lb) Point Molar Load



omit

**ENERGY AND ENVIRONMENT PART 1:
ENVIRONMENTAL TECHNOLOGIES**

PRECEDING PAGE BLANK NOT FILMED

97.

~~96~~ INTENTIONALLY BLANK



59-46
150479
N93-25570
p-10

Development of a Cone Penetrometer for Measuring Spectral Characteristics of Soils In Situ

Landris T. Lee, Jr., Philip G. Malone
and Stafford S. Cooper
USAE Waterways Experiment Station
Vicksburg, MS 39180

ABSTRACT

A patent was recently granted to the U.S. Army for an adaptation of a soil cone penetrometer that can be used to measure the spectral characteristics (fluorescence or reflectance) of soils adjacent to the penetrometer rod. The system can use a variety of light sources and spectral analytical equipment. A laser-induced fluorescence measuring system has proven to be of immediate use in mapping the distribution of oil contaminated soil at waste disposal and oil storage areas. The fiber optic adaptation coupled with a cone penetrometer permits optical characteristics of the in-situ soil to be measured rapidly, safely, and inexpensively. The fiber optic cone penetrometer can be used to gather spectral data to a depth of approximately 25 to 30 m even in dense sands or stiff clays and can investigate 300 m of soil per day. Typical detection limits for oil contamination in sand is on the order of several hundred parts per million.

INTRODUCTION

Cone penetrometers have been used in soils investigations for foundations and roadways for over fifty years. A typical geotechnical cone penetrometer consists of a hollow, instrumented, steel rod that is forced into the ground at a constant rate by employing hydraulic rams and a large reaction mass. The rod and conical tip are generally instrumented to measure the force the soil generates on the standard conical tip, and the force the frictional resistance of the soil produces on the side wall of the rod. Additional adaptations have permitted the measurement of the soil electrical resistivity and the pore pressure of fluids in the soil. The purpose of the present paper is to discuss new adaptations that allow the cone penetrometer to be used to measure spectral properties of soils in-situ.

An engineering cone penetrometer of modern design typically consists of a 20- to 30-ton truck equipped with all-wheel drive. Hydraulic jacks are used to lift the truck up from the ground so that all of the weight of the truck can be mobilized as a reaction mass. The hydraulic rams mounted in the truck use the reaction mass to force the penetrometer rod into the underlying soil. The electronics and computer equipment needed to readout and record data from the instruments in the rod is mounted in the van body that houses the rams (Figure 1).

The cone penetrometer is recognized in the geotechnical community as a rapid method for gaining access to the subsurface soils in order to make in-situ measurements or to recover samples of soil or groundwater. A typical cone investigation uses a 35-mm-diameter, hollow, steel rod that is forced into the soil at 2 cm/sec. A truck can generally investigate 300 m of soil in a working day. Using a 200-kN thrust the penetrometer can reach depths of at least 25 to 30 m even in dense sands and stiff clays (1).

The cone penetrometer equipped with suitable sensors is finding new applications in reconnaissance-level site investigations where contaminated soil and shallow groundwater are suspected. Cone penetrometers represent a faster, safer, and more economical alternative to drilling, sampling and sample analysis (2). Cone penetrometers produce no cuttings for disposal and a relatively simple adaptation to the ram unit allows the penetrometer rods to be cleaned as they are brought out of the ground. The penetrometer van can be adapted to maintain a cool, controlled air supply for the equipment operators. Air quality monitors located in the rod handling area assure that the van interior is a safe working area.

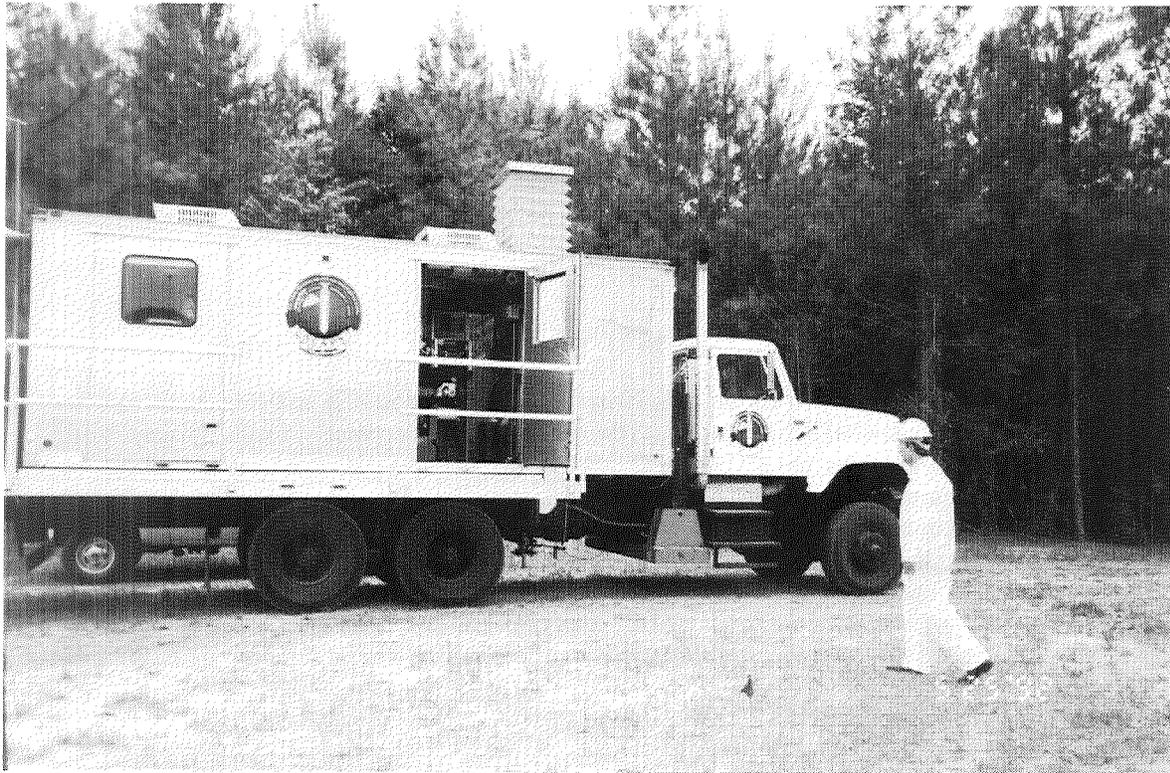


Figure 1. Photo of a penitrometer truck equipped for waste site investigations. The forward compartment houses the hydraulic rams for forcing the penitrometer rods into the ground. The rear compartment houses the instrumentation.

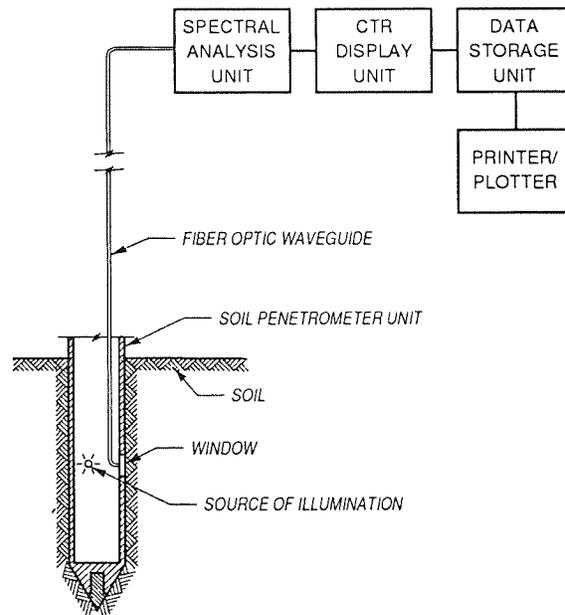


Figure 2. Schematic drawing of the system used for measuring soil spectral characteristics (after Figure 1 in U.S. Patent No. 5,128,882)

THE FIBER OPTIC CONE PENETROMETER

The U. S. Army Toxic and Hazardous Materials Agency is directing a tri-services research effort on applications for the cone penetrometer in site characterization and has tasked the U.S. Army Engineer Waterways Experiment Station to develop the cone penetrometer equipment that will allow it to be used as a screening tool for locating contaminated soil and groundwater. As one result of research in this area, the U.S. Army has been granted a patent on a novel method for examining the spectral characteristics of soil adjacent to the penetrometer rod (3). The patent describes a method for using a window that passes through the wall of the penetrometer tube. Light from inside the penetrometer tube is used to illuminate the soil opposite the window. The light returning from the soil is captured by a fiber optic waveguide inside the penetrometer and transferred to the surface. Spectral analysis equipment attached to the end of the fiber at the surface is used to determine the energy distribution of light returning from the soil. A schematic of the system is shown in Figure 2. The spectra are displayed in near real time in the penetrometer instrument compartment and are evaluated and recorded.

The illumination source and fiber optic waveguide in the cone penetrometer can be configured in a variety of ways to measure different phenomena. The basic fiber optic system can be used as a fluorometer or as a reflectometer. Two experimental fluorometer units have been built by the Army to detect contamination from hydrocarbons in soil, one unit employed miniature UV lamps housed in the penetrometer as an excitation source and used a grating spectrophotometer as a spectral analysis unit. A single waveguide was used to carry the fluorescent signal to the surface. A second design developed in an Army and Navy cooperative effort used a pulsed nitrogen laser emitting at 337 nm coupled to a fiber as an excitation source and a fixed grating with a charge-coupled photodiode as a spectral analyzer. Both single- and double-fiber designs have been built and evaluated. The Air Force has assembled a fiber optic cone penetrometer to detect jet fuel contamination in soil. The Air Force unit uses a portable, tunable dye laser coupled to a fiber bundle and a grating spectrophotometer as a spectral analyzer.

A prototype reflectometer has been assembled by the Army to evaluate the use of reflectometry in the visible portion of the spectrum as an aid to detecting wastes (such as TNT washout and rinse waters) that have a distinctive color. The reflectometer uses a tungsten lamp coupled to a fiber waveguide to provide illumination at the window and a second fiber waveguide and a spectrophotometer as a spectral analyzer.

Because of the widespread problem of soil contamination from fuel spills or leaks and waste oil disposal, the fluorometer is currently the configuration that has the widest application. The most useful system to date is the laser-induced fluorescence (LIF) unit that uses a nitrogen laser as a UV excitation source. The components used are specified in Table 1, and illustrated in Figure 3.

The penetrometer tool that is used in the LIF system (Figure 4) is adapted from a standard penetrometer cone that would be used to measure soil strength (4). The window through the wall of the penetrometer rod and the fiber optic elements are contained in a module that rides above the standard cone. The only modification that was produced in the lower part of the penetrometer is the addition of a grouting system that allows the tip of the penetrometer cone to be ejected and grout to be pumped down through the rod to seal the hole as the penetrometer is withdrawn. With the exception of the grouting system the penetrometer tool follows the standard design with regard to the tip configuration and the area of the lower part of the rod (the sleeve) that is used for soil friction measurements.

TABLE 1

LASER INDUCED FLUORESCENCE (LIF) UNIT COMPONENTS

Illumination Source: Laser Photonics LN 1000 Nitrogen Laser (337.1 nm excitation).

Fiber Optic Waveguide: Ensign - Bickford 360 micron core fiber optic, 400 micron total diameter.

Spectral Analysis: EG&G PARC Model 1460 Optical Multichannel Analyzer; Model 1302 Fast Pulsar; Model 1229 Spectrograph; Model 1421 Photodiode Array.

Data Acquisition and Processing: Hewlett Packard Vectra 486 Computers (2) networked via Ethernet; Data Translation A/D and D/A boards; Hewlett Packard LaserJet III printer.

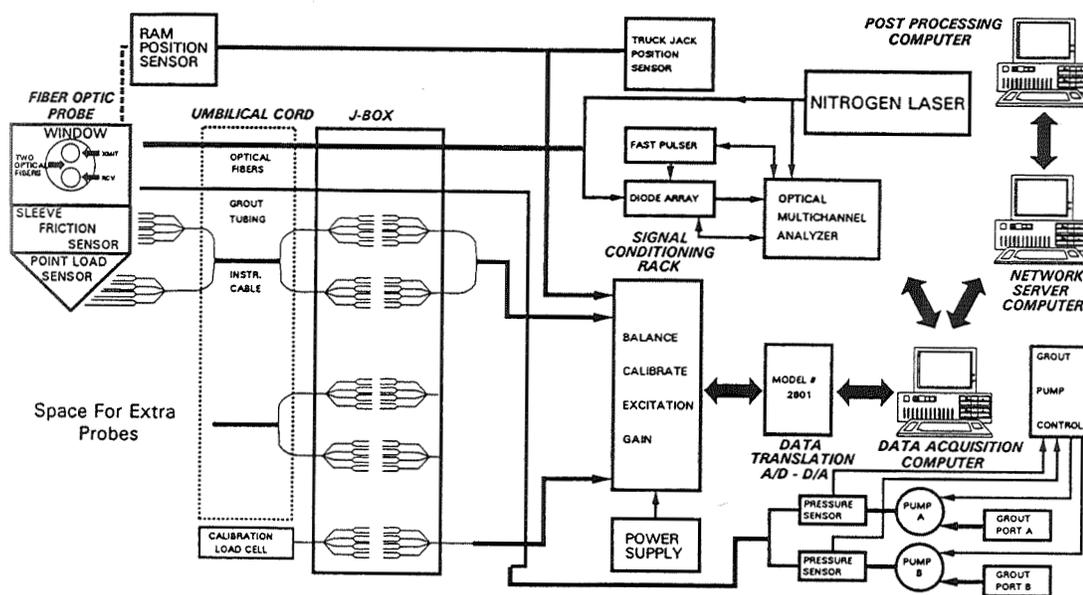


Figure 3. Layout of the instrumentation used in the laser-induced fluorescence (LIF) system.

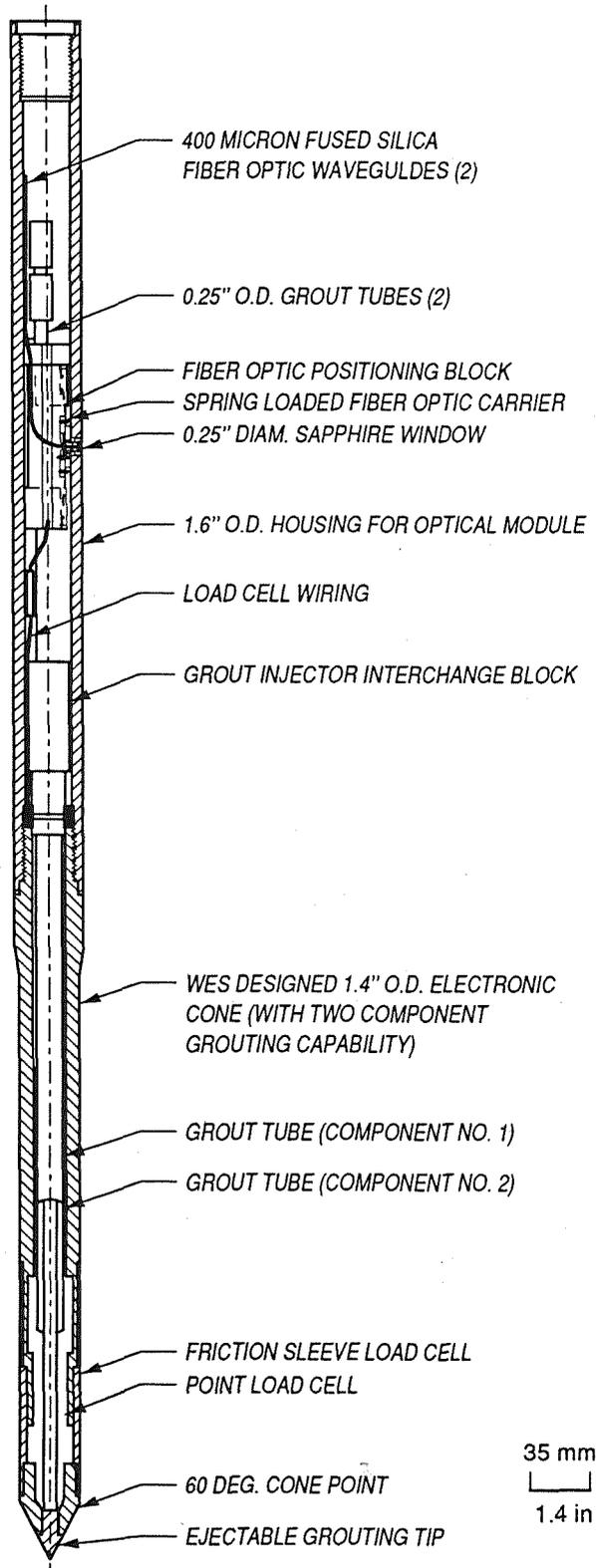


Figure 4. Cross section of the penetrometer tool showing the instrumented cone and the optical module.

The optical module (Figure 5) is designed with a replaceable sapphire window. The module consists of a fiber optic carrier that slides inside the penetrometer rod. The carrier holds two 400-micron diameter silica-on-silica fibers. The fibers are epoxied into a holder that assures that the cones of acceptance of the fibers overlap at the outer surface of the window. As the window is screwed into the carrier, the threads on the window pull the spring-loaded holder and the attached fibers into position behind the window. The wiring harness for the strain gages at the tip and sleeve of the penetrometer rod passes through slots in the optical fiber carrier. The optical fiber carrier has two stainless steel tubes inserted through it to allow the grout to be pumped down to the ejectable tip.

OPERATION OF THE FIBER OPTIC PENETROMETER

As in any technique that depends on spectral data for detection or quantification of a compound in a complex matrix, the detection limit depends on the response of the compound of interest and the influence of the matrix. The detection limit that can be obtained from the fluorometer when it is used to examine fuel-contaminated soil depends on the fluorophore present in the fuel (for example, polynuclear aromatic compounds in diesel fuel) and the conditions in the soil. With carefully prepared standards and heavy-grade fuels (rich in polynuclear aromatic compounds) in a sand matrix, detection limits of a few parts per million are possible (Figure 6). At many sites where free fuel is present in the soil as a non-aqueous phase liquid, low detection limits are not necessary if the objective is to determine the shape, size and depth of the mass of oil-saturated soil. The fluorometer can be used to find critical locations where the penetrometer or drilling techniques can then be used to collect soil or groundwater samples for analysis.

POTENTIAL APPLICATIONS FOR THE FIBER OPTIC PENETROMETER

The first application for the fiber optic penetrometer has been fluorometry for the detection of fuel in soil. Fuel residues at five sites have been mapped. By using the large volume of sensor data from the cone penetrometer and a volume mapping computer routine, a model of the contaminated soil mass can be prepared that shows the probable concentration of the fuel in the soil and the location and depth of each sensor reading. The visualization of the sensor-derived concentrations can be used in planning monitoring or remedial actions. An example of a three-dimensional plume map produced from penetrometer sensor data is shown in Figure 7.

The fiber optic cone penetrometer even in the form of a simple fluorometer offers potential for commercial applications. The U. S. Environmental Protection Agency estimates indicate that there are over two million fuel storage tanks in the United States. Surveys indicate that on an average one tank in three is leaking (5). This problem alone would justify commercializing the fiber optic cone penetrometer. The penetrometer also has uses in tracking the flow of landfill leachate and septic tank effluent. Fluorescent dyes have been used with the existing fluorometer to determine the direction and velocity of groundwater movement under a dredged material disposal area.

The fiber optic cone penetrometer will potentially become even more useful in the future as waveguides are developed that will allow broad range remote spectral studies, especially in the infrared. Development of more specific techniques using infrared, visible, and UV reflectometry can potentially expand the technological benefits. Experimental work is also underway in evaluating the use of resonance Raman spectrometry in a fiber optic cone penetrometer. Future fiber optic penetrometer applications will build upon basic designs developed during the production of the fluorometer system for the cone penetrometer (6, 7).

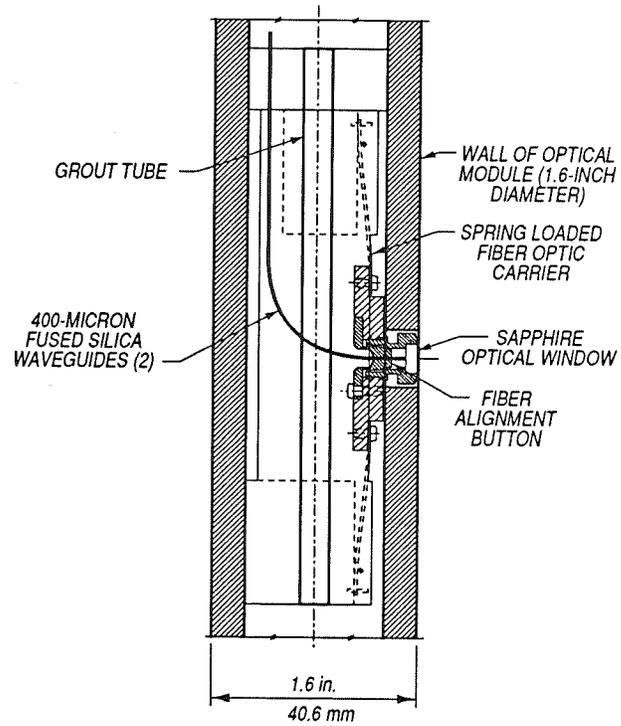


Figure 5. Detail of the fiber optic module showing the sapphire window. Note that the fibers are pulled in position as the window is screwed into place.

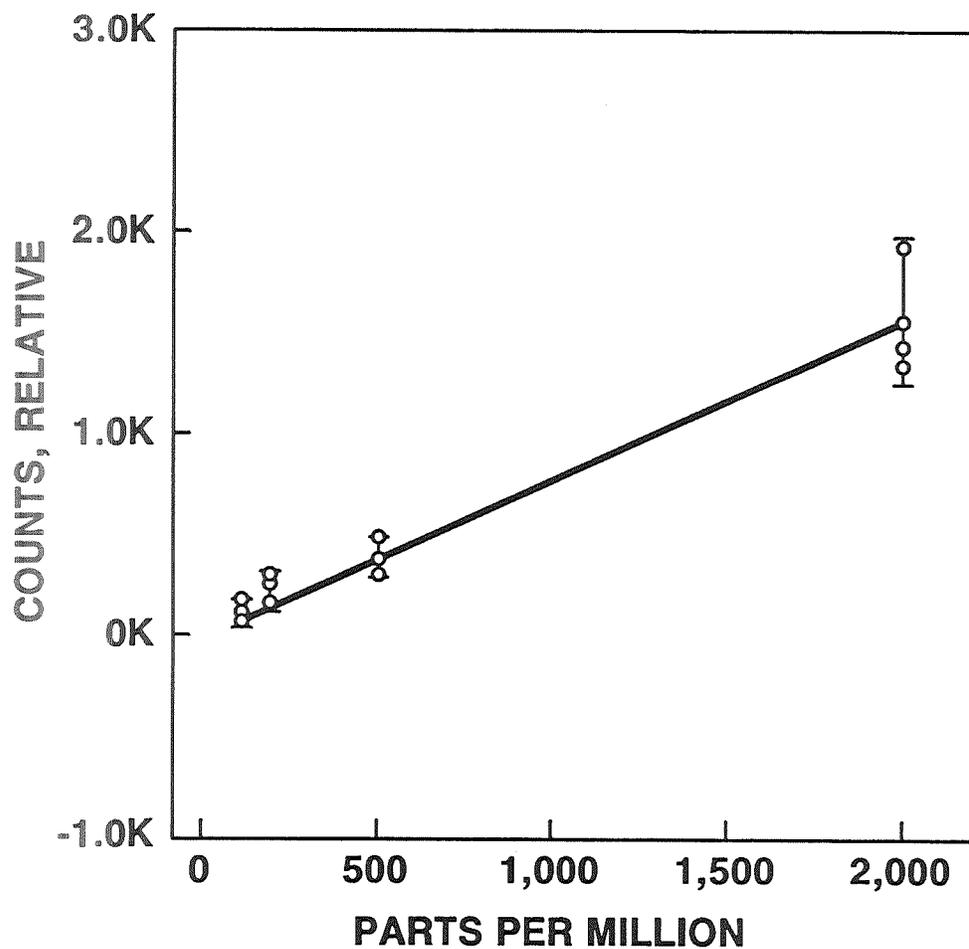


Figure 6. Typical calibration curve for fluorescence of #2 diesel fuel oil in sand. Parts per million are presented on a weight/weight basis. Error bars are +/- one standard deviation from the mean.

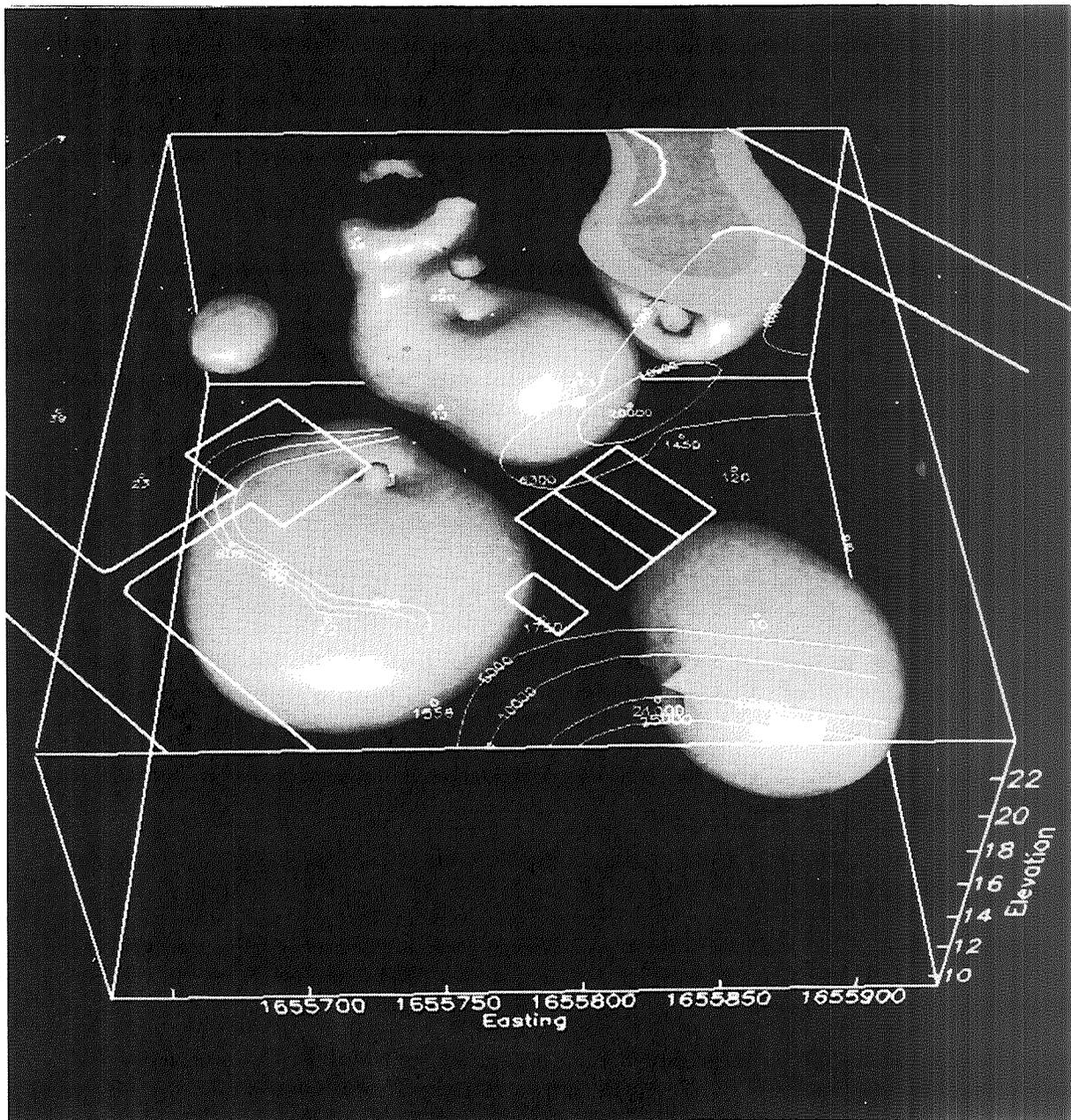


Figure 7. Block view of the masses of fuel-contaminated soil in a fuel storage area. The irregularly shaped masses are the volume of soil that fluoresced at levels equivalent to sand contaminated with over 2000 ppm of diesel fuel. The contours indicated on the surface show the concentration of hydrocarbon (measured as ppm hexane) in soil gas.

ACKNOWLEDGEMENTS

The development of the penetrometer for hazardous waste site characterization is sponsored jointly by the Department of Defense and the Department of Energy. The research effort is sponsored by the U.S. Army Toxic and Hazardous Materials Agency and is a tri-services project undertaken with the cooperation of the Air Force and Navy. This report incorporates data obtained in a number of field exercises conducted by technical staff at the USAE Waterways Experiment Station. The schematic in Figure 3 was prepared by Mr. Jeff Powell and Mr. Bobby Reed. The machine drawings in Figures 4 and 5 were supplied by Mr. Terry Kleinman. The sample calibration curve in Fig. 6 was prepared by Mr. Clifford Grey. The plume diagram in Figure 7 was prepared by Ms. Amy Chrestman.

The tests described and the resulting data presented herein, unless otherwise noted, were obtained from research conducted under the Site Characterization and Analysis Cone Penetrometer System Tri-Services Program directed by the U.S. Army Corps of Engineers. Permission to publish this report was given by the Chief of Engineers. The views of the authors do not purport to reflect the position of the U.S. Department of Army or of the U.S. Department of Defense. Citation of trade names does not constitute an official endorsement or approval of the use of any commercial products.

REFERENCES

1. de Ruiter, J. 1981. "Current Penetrometer Practice." in G. Norris and R. Holtz (eds.) Cone Penetration Testing and Experience: 1-49. New York, American Society of Civil Engineers, New York, N.Y.
2. Schroeder, J. D., S. R. Booth, L. K. Trocki. 1991. "Cost Effectiveness of the Site Characterization and Analysis Penetrometer System." Rept. No. LA-UR-91-4016. Los Alamos National Laboratory, Los Alamos, NM. 46 pp.
3. U.S. Patent Office. 1992. "Device for Measuring Reflectance and Fluorescence of In-situ Soil." Patent No. 5,128,882 issued July 7, 1992.
4. ASTM. 1991. Annual Book of ASTM Standards 04.08:439-444. American Society for Testing and Materials, Philadelphia, PA.
5. Valentinetti, R. A. 1989. "Federal Underground Storage Tank Regulations and Contaminated Soils." pp. 55-60 in Kostecki, P. T. and E. J. Calabrese. Petroleum Contaminated Soils, Vol. 2, Lewis Publishers, Chelsea, MI.
6. Lee, L. T. and others. 1992. "Developing Expanded Capabilities for the Site Characterization and Analysis Penetrometer System (SCAPS)." Prepared for publication in the Proceedings of the 16th Annual Army Environmental R&D Symposium, 23-25 June 1992, Williamsburg, VA, 12 pp.
7. Chrestman, A. M. and others. 1992. "Rapid Detection of Hydrocarbon Contamination in Ground Water and Soil." Prepared for publication in ASCE Waterforum '92, 2-5 August 1992, Baltimore, MD, 6 pp.

RECLAMATION WITH RECOVERY OF RADIONUCLIDES AND TOXIC METALS
FROM CONTAMINATED MATERIALS, SOILS, AND WASTES

A. J. Francis and C. J. Dodge

Department of Applied Science
Biosystems and Process Sciences Division
Brookhaven National Laboratory
Upton, NY 11973

N 9 3-25571

150480

P. 9

ABSTRACT

A process has been developed at Brookhaven National Laboratory (BNL) for the removal of metals and radionuclides from contaminated materials, soils, and waste sites. In this process, citric acid, a naturally occurring organic complexing agent, is used to extract metals such as Ba, Cd, Cr, Ni, Zn, and radionuclides Co, Sr, Th, and U from solid wastes by formation of water soluble, metal-citrate complexes. Citric acid forms different types of complexes with the transition metals and actinides, and may involve formation of a bidentate, tridentate, binuclear, or polynuclear complex species. The extract containing radionuclide/metal complex is then subjected to microbiological degradation followed by photochemical degradation under aerobic conditions. Several metal citrate complexes are biodegraded and the metals are recovered in a concentrated form with the bacterial biomass. Uranium forms binuclear complex with citric acid and is not biodegraded. The supernatant containing uranium citrate complex is separated and upon exposure to light, undergoes rapid degradation resulting in the formation of an insoluble, stable polymeric form of uranium. Uranium is recovered as a precipitate (polyuranate) in a concentrated form for recycling or for appropriate disposal. This treatment process, unlike others which use caustic reagents, does not create additional hazardous wastes for disposal and causes little damage to soil which can then be returned to normal use.

INTRODUCTION

The presence of radionuclides and toxic metals such as As, Be, Cd, Cr, Hg, Mn, Ni, and Zn in wastes, soils, and materials, at many Department of Energy (DOE) and other facilities is a major environmental concern. For decontamination of the waste material, both metal and radionuclide contaminants must be removed from the contaminated site so that the site can be returned to a useful condition. It would be desirable and beneficial to the environment to provide a comprehensive method for the removal of toxic metals and radionuclides from contaminated sites with reclamation of the soil.

Previous large-scale methods devised to deal with the problems of contaminated materials and soils have utilized caustic reagents such as hot sulfuric or hydrochloric acids, and oxidizing agents such as sodium hypochlorite, to extract the metals. While these methods can remove contaminants, they also cause irreparable damage to the soil, generating secondary waste streams which create additional hazardous waste disposal problems. For example, various soil washing methods were discussed at the DOE Soil Washing Workshop [1]. Madic et al. [2] used bidentate phosphamides enhanced by nitric acid to extract metal ions such as lanthanides, U(IV), Am(III) and Pu(IV). Kim et al. [3] immobilized radioactive Sr-90 by coprecipitation with Ca-, Al-, and Fe-phosphate in contaminated soils. Raghavan et al. [4] described three generic types of extractive treatments for cleaning excavated soils: water washing, augmented with a basic or surfactant agent to remove organics, and with acidic or chelating agents to remove organics and heavy metals, organic-solvent washing to remove hydrophobic organics and polychlorinated biphenyls; and air or stream stripping to remove volatile organics. Kochen and Navatil [5] removed americium and plutonium from contaminated soil by wet-screening, attrition scrubbing, wet-screening additives, and fixation by conversion to glass.

Other methods of metal removal from materials have also been reported. For example, U.S. Patent No. 4,973,201 describes a method for solubilizing precipitated alkaline earth-metal-sulfate scale in contaminated earth by contacting the earth with a polyaminocarboxylic acid chelating agent (EDTA, DTPA) and an oxalate ion synergist, and leaching the solubilized precipitate from the earth with water. To dispose of the dissolved sulfates, the leachate is either treated by chemical methods or returned to subterranean formation. Radioactive contaminated components of nuclear reactors have been decontaminated using citric acid (U.S. Patent Nos. 4,839,1000; 4,729,855; 4,460,500; 4,587,043; 4,537,666; 3,664,870 and 3,103,909). In these patents, metal recovery methods involve ion exchange columns, porous DC electrodes, or combusting the organics. Nishita et al. [6] used inorganic and organic compounds including citric acid to extract Pu from contaminated soils in order to correlate with Pu uptake by plants. Photochemical oxidation of uranium(IV) citrate by tungsten filament lamp, and the formation of an insoluble, stable polymeric form of uranium upon exposure of uranium citrate to light, has been reported [7,8].

Many of the chelating agents used in decontamination have been shown to undergo little degradation by microorganisms [9-11]. Biodegradation of these metal chelates should result in the precipitation of released ions as water-insoluble hydroxides, oxides, or salts, thereby retarding the migration of metals. Recently, we reported that the type of complex formed between the metal and citric acid plays an important role in determining its biodegradability [12]. The presence of free hydroxyl groups of citric acid is the key determinant in effecting biodegradation of the metal complex. For example, Ca, Fe(III), and Ni formed mononuclear bidentate complexes and were readily biodegraded; whereas, Cd, Cu, Fe(II), and Pb formed mononuclear tridentate complexes, and U formed a binuclear complex involving the hydroxyl group of the citrate, and were not biodegraded. The lack of degradation of tridentate and binuclear complexes was not due to toxicity, but probably limited by the transport and/or metabolism of the complex by the bacteria.

Various isolated concepts involving biodegradation, photodegradation, or chemical pretreatment have been reported in the literature. The problem of providing a thorough decontamination of the waste site as of yet has not been solved. In this paper, we present a total method for reclaiming radionuclide or toxic metal-contaminated materials, soils, sediments, and wastes with recovery of the contaminating metals to reduce toxic waste and with restoration of the soil [13].

TREATMENT PROCESS

The method for decontaminating radionuclides and other toxic metal-contaminated materials, soils, sediments, sludges and wastes, involves treating the contaminated material with citric acid and extracting the metals and radionuclides as citrate complexes. The solution is then treated with a *Pseudomonas fluorescens* ATCC No. 55241 and subjected to photolysis to degrade the complex and recover the radionuclides and metals through precipitation and biosorption reactions (Figure 1). The treated material is returned to its original use.

Extraction of Radionuclides and Metals

A sludge sample containing uranium and several toxic metals was obtained from a uranium processing facility. Ten grams of sludge was extracted with 100 ml of 0.40 M citric acid, for five hours in the dark using a wrist action shaker. The dry weight of the sludge before and after citric acid extraction was determined by drying at 60°C until a constant weight was obtained. The solids were digested in a mixture of hot nitric, perchloric and hydrofluoric acids in platinum crucibles. The citric acid extract and digested solids were analyzed for metals by ICP-MS. Table 1 shows extraction efficiency of various metals from sludge by citric acid. In this sample, metals Ag, As, Au, B, Bi, Cu, Gd, Hg, Li, Mo, Pb, and V were poorly extracted by citric acid treatment. Lack of extraction of these metals is probably due to the nature of mineralogical association with stable mineral phases in this particular waste [14]. For example, Cu was predominantly associated with the organic fraction and a small amount with the iron oxide and inert fractions and was not extracted by citric acid treatment.

Biodegradation of Metal Citrate Extract

Duplicate samples of undiluted (as is) and diluted (1 to 4) citric acid extract from the sludge were amended with nutrients consisting of 0.1% of NH_4Cl , K_2HPO_4 , and KH_2PO_4 . The pH was adjusted to 6.5 with NaOH and then the extracts (100 ml) were inoculated with 4 ml of 18-hour old culture of *Pseudomonas fluorescens* ATCC 55241. The samples were incubated on a shaker at 24°C. The bacterial inoculum was grown in medium containing citric acid, 2g; NH_4Cl , 1g; KH_2PO_4 , 1g; K_2HPO_4 , 1g; NaCl, 4g; MgSO_4 , 0.2g; distilled water, 1000 ml; and pH 6.5. All sample preparations were performed under low light to minimize any photochemical reactions. Periodically, 5 ml aliquots were removed, filtered through a 0.22 μm filter, and analyzed for (i) pH, (ii) citric acid biodegradation by HPLC using uv and refractive index detectors, and (iii) soluble uranium. At the end of the incubation period (after 118 hours for the diluted sample and after 322 hours for the undiluted sample) the supernatant and the solids consisting of bacterial biomass and any precipitated metals were separated from solution by centrifugation. The dry weight of the solids was determined and digested in a mixture of hot nitric and perchloric acids. The supernate and the digested solids were analyzed for uranium and other metals by ICP-MS (Table 2). The bacteria degraded citrate at a rate of 0.5-0.7 mM per day (Figure 2). The rate of degradation was much higher in the diluted extract than the undiluted sample. There was little change in concentration of uranium in samples subjected to biodegradation, indicating that the uranium citrate complex was not biodegraded (Figure 2). This is consistent with the previous results that uranium and certain metal complexes of citric acid are resistant to biodegradation [12]. About 9% of the total uranium was present in the bacterial biomass. Increased levels of cobalt, nickel, zinc, and zirconium, in the biomass digest, indicated that their citrate complexes were readily biodegraded.

Photodegradation of Uranium Citrate Extract

The supernate from the biodegradation treatment containing primarily uranium citrate complex was exposed to light to degrade the complex and recover uranium (Figure 3). The pH of the supernate was adjusted to 3.5 with HCl and the sample exposed to seven 60 watt, high output fluorescent growth lights. Periodically 2 ml samples were withdrawn, filtered through a 0.22 μm filter and analyzed for uranium, citric acid, and photodegradation products. At the end of the experiment (after 157-hours of exposure to light) the solutions were filtered through a 0.22 μm Millipore filter and analyzed for citric acid degradation products by HPLC, and for metals by ICP-MS. The uranium precipitated out of solution as a polymer soon after it was exposed to light (Figure 3). After 50 hours, ~ 85% of the uranium was removed from solution. In Table II, the removal efficiency of various metals from the citric acid sludge extract first subjected to biodegradation, followed by photodegradation, is presented.

Weight Loss

The solids remaining after extraction with citric acid were washed with deionized water, transferred to weighing dishes, and dried in an oven overnight at 105°C to determine the weight loss due to citric acid extraction. The extraction of the wastes with citric acid resulted in significant reduction in weight. Almost half (47%) of the sludge showed loss in weight due to solubilization and removal of toxic and nontoxic bulk components in the waste.

SUMMARY

These results show that (i) uranium was extracted from the waste sample with >85% efficiency using 0.4 M citric acid; (ii) other metals such as chromium, cobalt, manganese, nickel, strontium, thorium, zinc, and zirconium were also extracted from the waste; (iii) the uncomplexed excess citrate and several metal citrate complexes (Co, Ni, Zn and, Zn) with the exception of binuclear complexes were readily biodegraded by the bacterium, *P. fluorescens* ATCC 55241 and were recovered with the bacterial biomass; and (iv) the uranium citrate complex was photodegraded, allowing the uranium to form a polymer which was recovered as a concentrated solid.

This process has significant potential for commercialization because (i) it can be applied to a variety of materials and waste forms, (ii) it does not generate secondary waste streams, (iii) it causes little damage to soil, and (iv) environmentally and economically important metals are removed in a concentrated form for recovery and recycling. The use of combined chemical, photochemical, and microbiological treatment processes of contaminated materials will be more efficient and result in considerable savings in clean-up and disposal costs.

REFERENCES

1. Francis, C. W. An assessment of soil washing to remove uranium and mercury from Oak Ridge soils. Department of Energy's Soil Washing Workshop, Aug. 1990.
2. Madic, C., Rubin, P., Rodehueser, L., Delpuech, J. J., BenNasr, C., Fasan, G., and Bokolo, K. Extraction of metal ions by neutral B-Diphosoramides, Energy Research Abstracts 3636, Feb. 1991.
3. Kim, K. H., Ammons, J. T., and Lee, S. Y. Immobilization of radioactive strontium in contaminated soils by phosphate treatment. Energy Research Abstracts 1993, Jan. 1991.
4. Raghavan, R., Coles, E., and Dietz, D. Cleaning excavated soil using extraction agents: A state-of-the-art review. Final Report June 1985 - January 1989, Energy Research Abstracts 5106, Jan. 1990.
5. Kochen, R. L. and Navatil, J. D. Americium and plutonium removal from contaminated soil. Energy Research Abstract 36247, Sep. 1985.
6. Nishita, H., Havg, R. M., and Rutherford, T. Effect of inorganic compounds on the extractability of ^{239}Pu from an artificially contaminated soil. J. Environ. Qual. 6, 451-55 (1977).
7. Adams, A. and Smith T. D. The formation and photochemical oxidation of uranium(IV) citrate complexes. J. Chem. Soc. 4, 4846-50 (1960)
8. Ohyoshi, A. and Ueno, K. Studies on actinide elements(IV): Photochemical reduction of uranyl ion in citric acid solution. J. Nucl. Inorg. Chem. 36, 379-84 (1974).
9. Francis, A. J. Microbial transformation of low-level radioactive wastes in subsoils. In Soil Reclamation Processes: Microbiological Analyses and Applications, R. L. Tate and D. Klein, Editors, pp. 279-331, Marcel Dekker, New York, 1985.
10. Francis, A. J. Microbial transformation of toxic metals and radionuclides in mixed wastes. Experientia 46(8), 840-851 (1990).
11. Francis, A. J., Dodge, C. J., Gillow, J. B., and Cline, J. E. Microbial transformation of uranium in wastes. Radiochimica Acta. 52-53, 311-16 (1991).
12. Francis, A. J., Dodge, C. J., and Gillow, J. B. Biodegradation of metal citrate complexes and the implications for toxic metal mobility. Nature 356, 140-2 (1992).
13. Francis, A. J. and Dodge, C. J. Waste site reclamation with recovery of radionuclides and metals. Patent Pending, AUI 91-18 (1992).
14. Francis, A. J., Dodge, C. J., and Gillow, J. B. Microbial stabilization and mass reduction of wastes containing radionuclides and toxic metals. U.S. Patent No. 5,047,152 (1991).

ACKNOWLEDGEMENTS

We thank C. C. Neill for assistance. This research was performed under the auspices of the Environmental Sciences Division's Subsurface Science Program, Office of Health and Environmental Research, Office of Energy Research, U. S. Department of Energy, under Contract No. DE-AC02-76CH00016.

TABLE 1

Extraction Efficiency of Metals From Sludge by Citric Acid

Metal	Total Metal in sludge ($\mu\text{g}\cdot\text{gdw}^{-1}$)	% Metal Extracted
Ag	41 \pm 30	2.4
Al	30500 \pm 500	58.7
Au	1800 \pm 500	<1
Ba	427 \pm 25	24.4
Be	5.21 \pm 0.45	60.1
Cd	66 \pm 6	9.1
Co	10.7 \pm 0.3	74.8
Cr	342 \pm 10	74.6
Cu	329 \pm 18	1.2
Ga	28.8 \pm 0.6	25.7
Mg	7510 \pm 100	89.2
Mn	234 \pm 3	82.9
Ni	1120 \pm 10	80.0
Pb	224 \pm 27	<1
Pd	5.51 \pm 0.70	49.0
Sb	5.67 \pm 0.05	68.8
Sn	17.6 \pm 0.4	93.1
Sr	125 \pm 5	59.2
Th	3.08 \pm 0.10	94.2
Ti	922 \pm 95	28.4
U	2410 \pm 100	86.8
V	121 \pm 7	4.1
Zn	839 \pm 7	59.6
Zr	209 \pm 4	84.2

Sludge extracted for 5 hours with 0.4 M citric acid solution.

TABLE 2

Effects of Biodegradation Followed by Photodegradation
in the Treatment of Citric Acid Sludge Extract

Metal	Before Treatment ¹ (μM)	% Removal		Total
		After Biodegradation ²	After Photodegradation ³	
Al	7410	92	1	93
Ba	12.0	92	< 1	92
Be	3.22	>99	< 1	>99
Co	0.866	71	6	77
Cr	60.8	< 1	25	25
Ga	0.889	73	16	89
Mn	37.2	98	< 1	98
Ni	192	64	1	65
Pd	0.311	64	30	94
Sb	0.361	2	14	16
Sn	1.71	>99	< 1	>99
Sr	10.2	98	1	99
Th	0.112	96	3	>99
Ti	80.2	96	< 1	96
U	94.8	9	78	87
Zn	86.1	90	5	95
Zr	61.7	97	2	99

¹ Sludge was extracted for five hours with 0.4M citric acid.

² Samples analyzed 118 hours after inoculation with Pseudomonas fluorescens ATCC No. 55241 but before photodegradation.

³ Samples analyzed after biodegradation and 157 hours of exposure to light.

ND-none detected

Waste Site Reclamation With Recovery Of Radionuclides and Metals

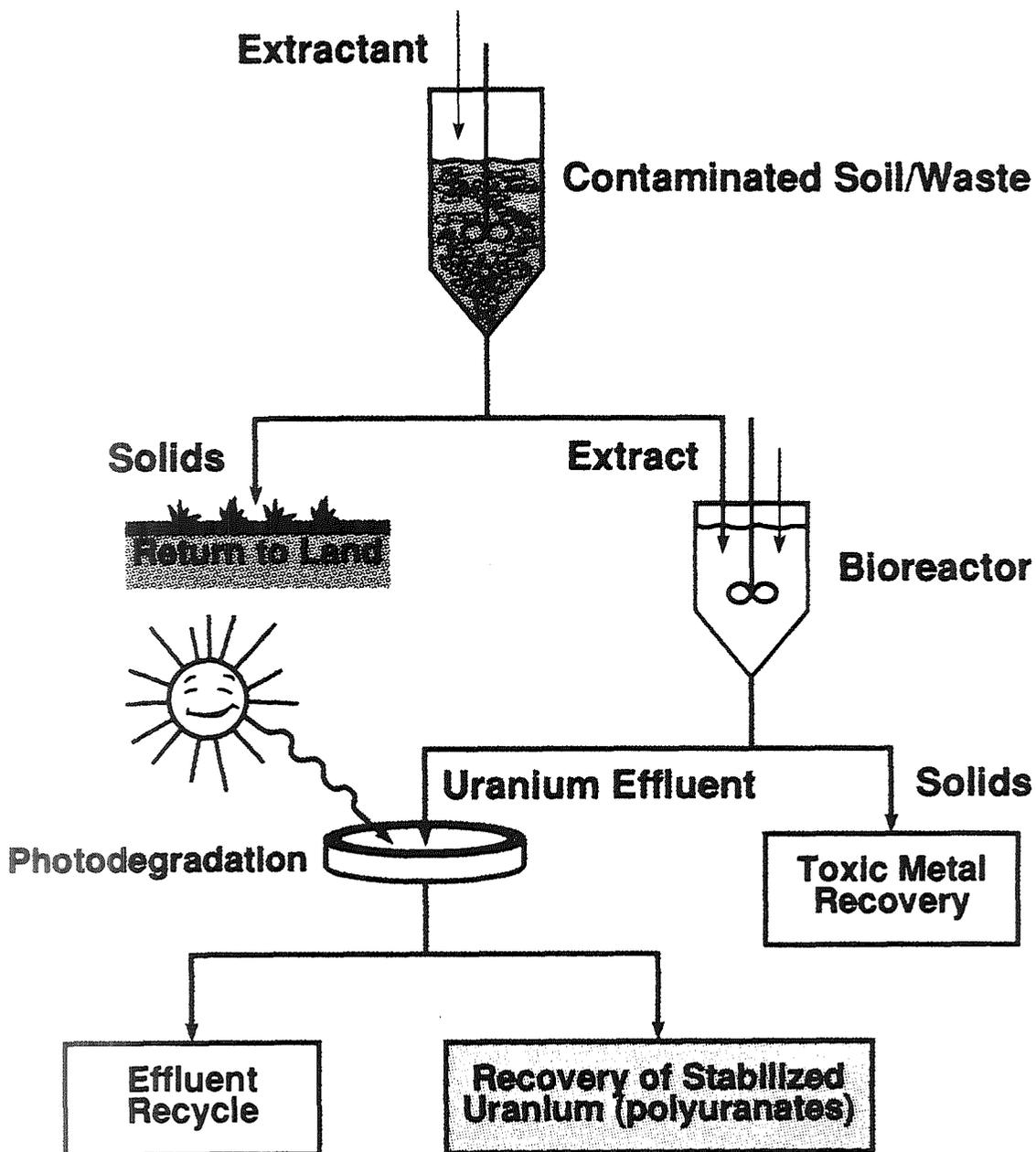


Figure 1. Schematic of the Treatment Process

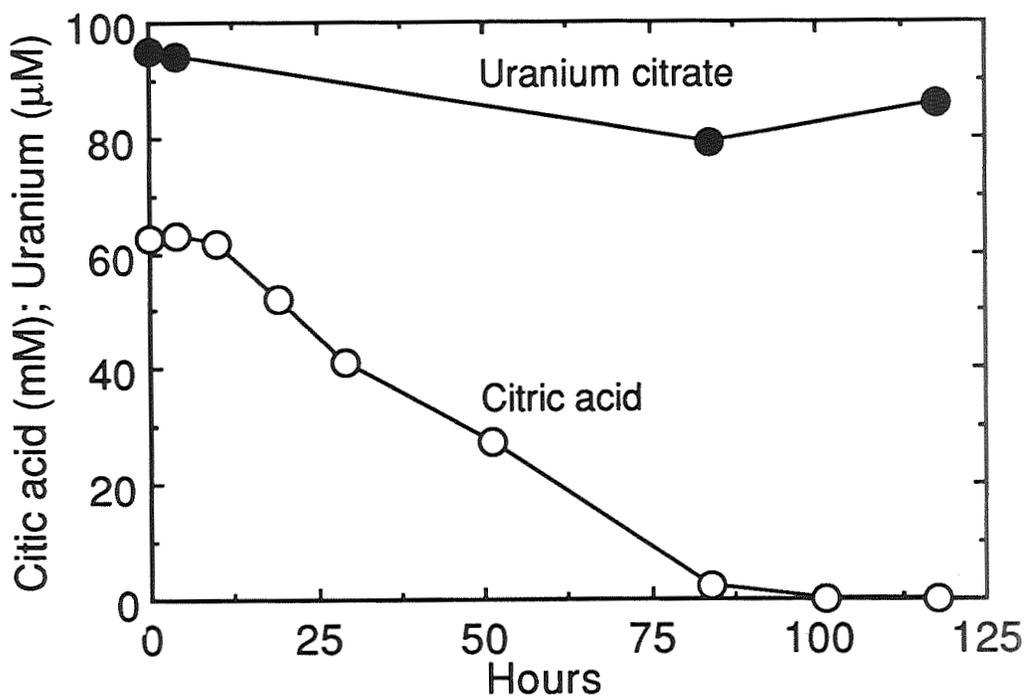


Figure 2. Biodegradation of citrate sludge extract

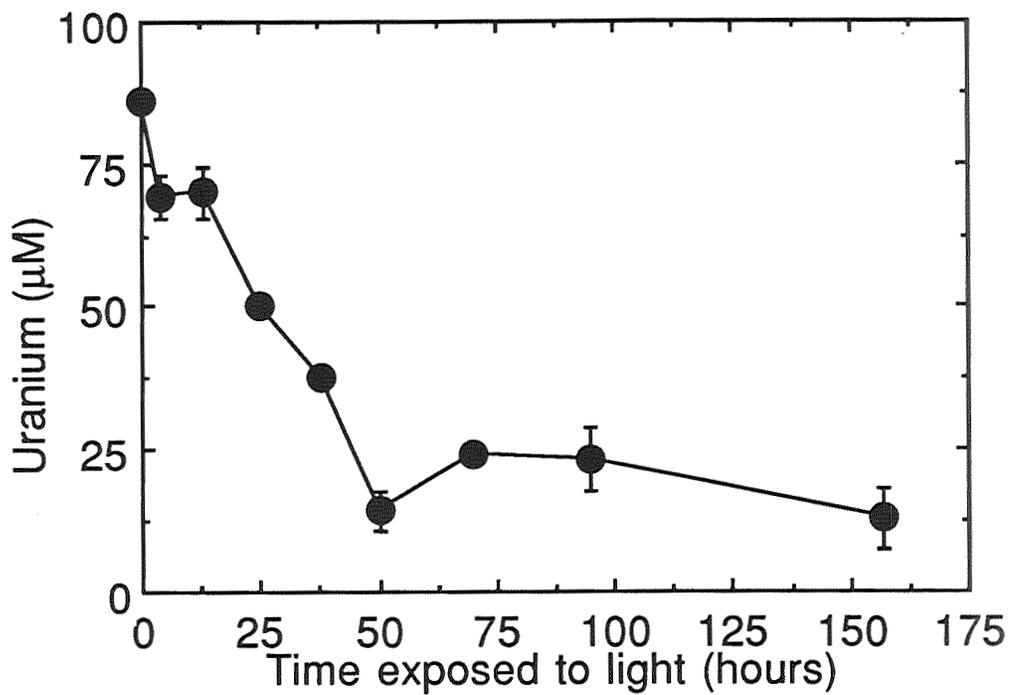


Figure 3. Photodegradation of uranium citrate sludge extract after biodegradation

511-44

150481

N 93-25572

P-3

CONVERTING ENVIRONMENTAL WASTES INTO VALUABLE RESOURCES

Leonard A. Duval
Recotech Corporation
Aurora, OH 44202

ABSTRACT

This concept employs a viable energy saving method that uses a solvent to separate oil from particle matter; it can be used in metal forming industries to deoil sludges, oxides, and particle matter that is presently committed to landfill. If oily particles are used in their oily state, severe consequences to environmental control systems such as explosions or filter blinding, occur in the air handling equipment. This is due to the presence of hydrocarbons in the stack gasses resulting from the oily particles.

After deoiling, the particles can be recycled and the separated oil can be used as a fuel.

The process does not produce a waste of it's own and does not harm air or water. It demonstrates the dual benefits of being commercially viable and in the national interest of conserving resources.

INTRODUCTION

This presentation covers a method to recover valuable resources from oily wastes by means of solvent extraction of the oil from the wastes. This method is known as the Duval Process and allows for:

- Recovery of solids
- Recovery of oil
- Elimination of harmful effects of oil on solids
- Elimination of harmful environmental problems caused by oil in disposal sites

This process has been licensed to major engineering and construction firms.

Some practical applications are as follows:

- Removal of oil from sludges and fine particle matter
- Decontamination of soils when the contaminant is a hydrocarbon
- Removal of solids from oily material

The process is one that maximizes the use of wastes that, at great cost, must be landfilled or destroyed to safeguard the environment. The process employed is an economical approach to maximizing resource recovery from oily wastes by use of solvent extraction of the oil from the solid particles.

Large generators of oily wastes are steel mills and machine shops. Accidental spills of oily materials have caused contamination of soils that require cleansing as insurance against future environmental damage.

In order to change the size or shape of metal components, machining operations such as shaving, forging, pressing, grinding, or polishing are required. A lubricant, usually an oil, is used to assist the tools during the machining operation by cooling and lubricating the tools and materials. Small particles of the metal are removed with the lubricant and later separated by gravity-settling so that the lubricant can be reused. Solids from a metal forming operation are coated with substantial quantities of the oil used in a metal

forming operation. Larger particles may contain as much as 2% of their weight in oil while smaller sizes may contain as much as 25% of their weight in oil and form a sludge like mass.

Quantities of metal fractions generated at a location are directly related to the type of forming and parts being manufactured; ie., in the manufacture of aircraft hydraulic fittings the metal removed may be as much as 50% of the original metal weight while metal removed during a grinding operational may only be a small fraction of the metal part. The oil quantity on these particles would likewise vary from 2-4% on machinings to 25% on grindings.

Much of the oily particle matter is being wasted due to difficulties resulting from its reuse, and this waste is sufficiently harmful to the environment to draw attention of environmental groups, research institutes, engineering companies, waste treatment firms, and firms engaged in resource recovery.

Thousands of tons of valuable resources are being sent to landfill disposals by the metal forming industries.
and there are numerous reports of landfills leaching contaminants into streams and the water table.

The cost of handling and managing the landfills has increased dramatically and therefore costs of disposing wastes to landfills is substantial. In addition to the direct cost of landfilling is the contingent legal liability of using the landfill. There have been cases where small waste producers have contaminated a non-hazardous landfill by disposing of hazardous materials which have contaminated the landfill and thus made all disposers that used the landfill responsible for the condition. The legal liabilities associated with public landfill disposal are so great that larger firms often find it necessary to operate their own landfills.

At present, operators of machining plants will use various methods to reduce the oil on the recyclable scrap, by storing the scrap in heated rooms and centrifuging. This is to reduce the loss of expensive cutting oil on the scrap as well as to make the scrap acceptable to the smelter. In addition, the waste generator is penalized by the smelter due to the pollution caused by the oil on the scrap shipped to the smelter.

If the oily wastes are used, serious damage to exhaust gas handling systems will result. Fires in electrostatic filters and blinding of bag filters in standard systems are dangers when the consuming facility involves high temperatures as in smelting or incineration. The oil on wastes being smelted cause detrimental metallurgical effects because of the sulfur in the oil. When incineration of oil at the smelter is required, the elevated temperature causes decomposition of particles requiring a replacement of lost components. Unavoidable wastes result from metal forming due to the oil on the wastes being generated.

If the wastes are piled or landfilled, the dedication of a large area for the landfill, costs associated with managing the landfill, hazards to ground water and ever changing regulation become prohibitive.

The solution to these problems is a program to make the wastes recyclable. The key to maximizing the recycling of the oily waste is the elimination of the oil on the waste being generated so that it can be either recycled directly back into the generator's production, or used as a raw material for another industry. Steel manufacturing produces a large amount of oily iron oxides that can be recycled back into steel making as a replacement for raw iron ore. Steel making oxides are also used as mason hardeners and colorings, in mixes for exothermic welding, as dense media and also as iron supplements in animal feed.

DESCRIPTION OF THE PROCESS

The process consists of a low-temperature, multi-stage (usually 3 or 4 stages) extraction using

methylene chloride flowing counter current to the oily solids, that extracts the oil and water linked to the solids. The solvents are settled and filtered for fines separation. The solids are dried by evaporating the residual solvent and water which are then condensed and recycled in the process.

The oily solvent is recovered by standard azeotropic distillation of the solvent and oil mixture carried out at 39 degrees Centigrade. The water is removed by gravity separation from the solvent.

Since the solvent has some affinity for the water, it is subjected to chilling in order to promote the gravity separation of the solvent and water. The water is further subjected to a polishing operation consisting of air-stripping the solvent to a carbon pack. The recovered water is used for make-up requirements within the process. This result is a reduced need for fresh water and a zero discharge of water into the environment.

The solvent also contains an amount of water which must be eliminated in order to insure against acidification of the solvent. The solvent containing the water is passed through a medium of recyclable desiccant which absorbs the water leaving the solvent relatively free of water. This stage of the process includes an automated system for maintaining solvent stability by addition of stabilizing agents as required.

Vents of the system pass through a granulated carbon column for collection of any solvent. The solvent is then recovered from the carbon column and returned to the process.

The entire system is close looped to maximize the solvent recovery. The circuits for recovering the solvent from the water and air, the value of the recovered oil as a fuel and elimination of landfill costs makes possible the conversion of an environmental problem into a profit center.

512-45
150482

LOW COST DEWATERING OF WASTE SLURRIES

J. B. Peterson, S. K. Sharma,
R. H. Church, and B. J. Scheiner

N 93 - 25 573

U.S. Bureau of Mines
Tuscaloosa Research Center
P.O. Box L, University of Alabama Campus
Tuscaloosa, AL 35486-9777

ABSTRACT

The U.S. Bureau of Mines has developed a technique for dewatering mineral waste slurries which utilizes polymer and a static screen. A variety of waste slurries from placer gold mines and crushed stone operations have been successfully treated using the system. Depending on the waste, a number of polymers have been used successfully with polymer costs ranging from \$0.05 to \$0.15 per 1,000 gal treated. The dewatering is accomplished using screens made from either ordinary window screen or wedge wire. The screens used are 8 ft wide and 8 ft long. The capacity of the screens varies from 3 to 7 gpm/ft². The water produced is acceptable for recycling to the plant or for discharge to the environment. For example, a fine grain dolomite waste slurry produced from a crushed stone operation was dewatered from a nominal 2.5 pct solids to greater than 50 pct solids using \$0.10 to \$0.15 worth of polymer per 1,000 gal of slurry. The resulting wastewater had a turbidity of less than 50 NTU and could be discharged or recycled. The paper describes field tests conducted using the polymer-screen dewatering system.

INTRODUCTION

In the processing of minerals to produce concentrates and products, often times a dilute slurry is generated that must be disposed of in some manner. The use of impoundments for waste disposal is common throughout the minerals industry. With the promulgation of new environmental regulations, the cost associated with using and maintaining impoundments for waste disposal is increasing dramatically. The Bureau of Mines has been conducting research to develop low cost dewatering techniques (1-2). The most recent research activities include effluents from placer mines and slurries generated in the production of crushed stone.

In placer mining, gold bearing gravels are treated in a washing plant to remove boulders, small rocks, sand, and fines while trapping the gold particles. This is usually accomplished by placing the gravel into a trommel or on vibrating screens where the gravel is sized from 0.5 to 1 inch. The undersized material is washed into a sluice box while the small rocks, sand, and fines flow off the end of the sluice box into a sump, where a majority of the rocks and sand settle out. The water containing the fines and some sand flows from the sump and into a pond system at the mine site. The settleable material drops out as the water moves sequentially through the system of ponds, leaving the fine grain silts and clays in suspension. This is commonly referred to as the non-settleable fraction of the gravel being treated. With time, more fines will settle resulting in a solution containing ultrafine or colloidal particles that will remain suspended indefinitely. Contamination of surface waters is possible if this turbid water is discharged. The extent of this problem depends on the character of the gravel being treated. For some gravels, very little colloidal particles will be formed, whereas, for others a significant amount can be generated.

The crushed stone industry produces a variety of sized stone for sale. The stone is mined, crushed, and sized (3). During the sizing operation, especially for the smaller fractions, the stone is washed to remove undesirable fines. This results in a fine grained slurry which must be impounded. Oftentimes, due to the location of the mine near municipalities, the land available for use as impoundments is limited requiring the impoundments to be emptied and the material transported to a location where it can be disposed. This often entails the use of equipment such as draglines to remove the material and trucks for transport of the consolidated material to the disposal site. The fine material is often thixotropic and causes difficulty in removal and transport (4). Impoundment design and maintenance and sediment disposal are regulated under a number of environmental and zoning regulations. The fines generated during the processing oftentimes contain certain fairly pure materials that are saleable if they can

be dried, i.e., waste fines from a magnesium/calcium carbonate quarry is often of high purity and can be sold as fine grain carbonate.

This paper describes a technique for dewatering fine grained slurries. The technique produces a consolidated material that can be handled by conventional equipment. The technique involves flocculation of the solids in a slurry with the proper polymer and dewatering on a static screen.

FLOCCULANT SCREENING

Prior to field testing, samples of the slurry are obtained from the mine site and laboratory experiments conducted. A large number of commercially available polymers are screened to determine the optimal polymer for the particular waste slurry. The technique used to test the polymer is described in detail elsewhere (4). The criteria used to determine the best polymer include dosage requirements, cost of polymer, clarity of discharge water produced, and the percent solids of the dewatered material. For the field testing described in this paper, high molecular weight polyethylene oxide (PEO) was the polymer chosen for the placer effluent tests, and a high molecular weight anionic polyacrylamide was used for the slurry from a crushed stone operation.

FIELD TESTING

The placer site evaluated is located in Livengood, AK, and the crushed stone operation is in Birmingham, AL. At the Alaska site, the unit was designed to handle up to 1,000 gpm of placer effluent. The flow sheet for the operation is shown in Figure 1. A 6-in pump with a capacity of 1,100 gpm was used to deliver the waste slurry to the unit. The PEO pumps with variable speed drive systems and with capacities of 45 and 20 gpm, respectively, were used to inject the PEO solution in-line by using 2-in-diam pipe. Placer effluent was delivered to the system by 6-in pipes of various lengths. A Kenics¹ static mixer (2-10 elements) was used, when needed, to increase the

¹Reference to specific products does not imply endorsement by the Bureau of Mines.

turbulence in the pipe. The flocculated slurry was emptied into a trough at the top of the screens and overflowed onto the screens. The water passed through the screen into a trough and was allowed to flow by gravity into the secondary pond. The dewatered solids rolled down the screen and discharged into a pit. The screen was comprised of two sections, the top section was set up at an angle of 47° and the bottom section at an angle of 38°. A common aluminum window screen, 16 by 18 mesh, was used as the screening device for the flocculated solids.

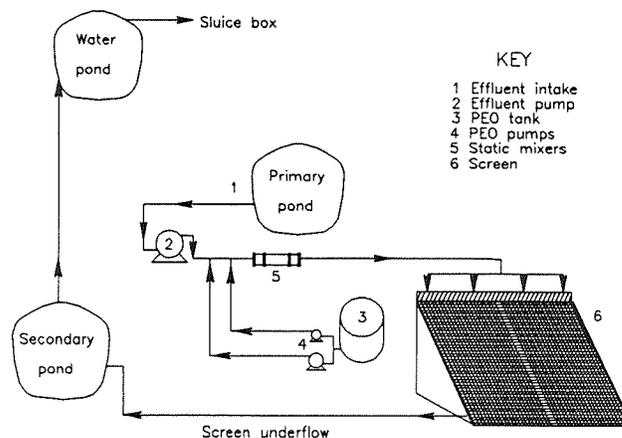


FIGURE 1.--Flow sheet for Alaska field test study.

For the crushed stone field test, a mobile unit was constructed and placed on a flat bed truck as shown in Figure 2. The equipment includes holding tanks for polymer solution, a mixing tank for preparing concentrated polymer solutions, a mixing tank for dilution of the polymer, and static dewatering screens 8- by 4-ft hinged together to form an 8- by 8-ft screen. The screens have horizontal openings that are 0.020 and 0.030 in wide, 2.75-in long. The screens can be replaced with other screens having horizontal openings of 0.010 or 0.040 in width. The waste slurry is pumped from the processing plant to the mixing trough above the screens. Polymer is added to this line. The polymer mixes with the waste and then enters the trough, and then overflows to the screen. The released water flows through the screen and is recycled back to the plant's water system. The solids roll off the bottom of the screen and are placed in a pit. Figure 3 shows the unit in operation at the crushed stone site.

RESULTS AND DISCUSSION

A variety of conditions were used at the Alaska field testing site. From previous research results it was known that obtaining the proper mixing of polymer and slurry was critical (5). To determine the best mixing conditions, a variety of different arrangements were tested. The PEO was injected into the feed line upstream of the dewatering screen. In-line mixing through 400 to 1,000 ft of pipe combined with a wide range of Kenics static mixers (2 to 10 elements) were tested. During initial testing on the primary pond, 600 ft of 6-in pipe and the 2-element static mixer produced best results. However, when the unit was transferred to the secondary pond, 600 ft of pipe alone produced good flocs. In both tests, treated water was recycled back to the secondary pond. The feed, which varied widely in solids content from 0.09 to 5.40 pct by weight and turbidity from 300 to 26,500 NTU, was dewatered using 0.01 pct PEO solution.

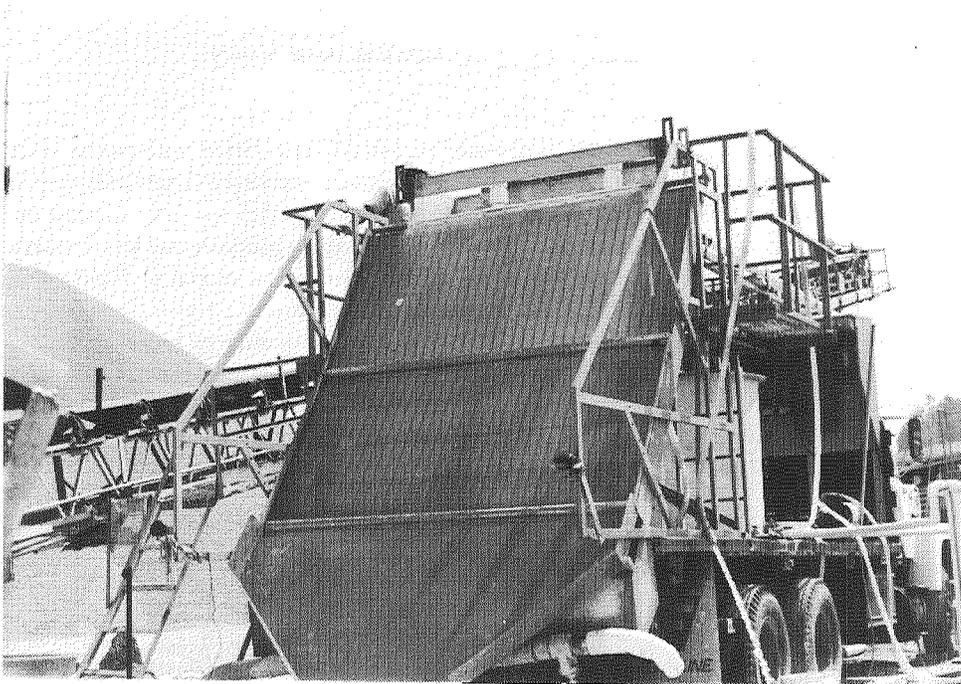


FIGURE 2.--Mobile dewatering unit at stone quarry.

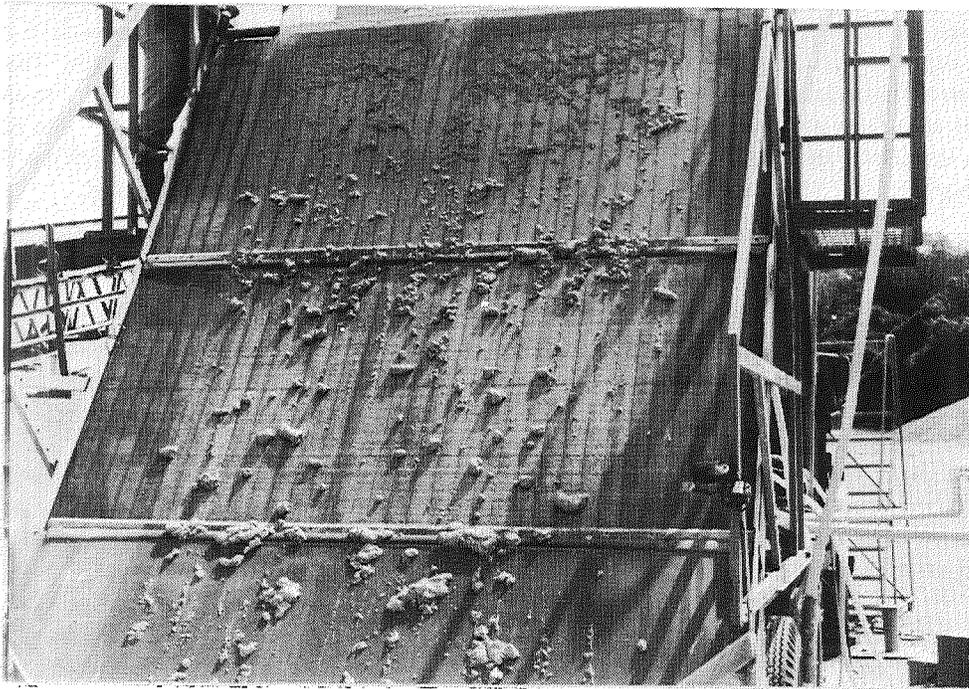


FIGURE 3.--Dewatering screen operating at stone quarry.

As shown in Table 1, the PEO dosage required to dewater placer effluent varied with initial solids and was calculated in pounds per 1,000 gallons of slurry treated. The PEO dosage increased from 0.02 to 0.19 lb/1,000 gal with an increase in initial solids. It was also found that the PEO dosage did not only depend on the initial percent solids but also depended on the Reynolds number (Re) which is directly proportional to the slurry flow rate. Therefore, Re is just a gauge of shear force requirement of a system. The results shown in Figure 4 indicate that when the Re increased from 50,000 to 130,000, the PEO dosage decreased from 0.15 to 0.04 lb/1,000 gal. Each point on Figure 4 represents an individual test and all the tests had a water quality of less than 50 NTU. Also, Table 1 shows that for higher initial solids in the 26,000 range the static mixer was a benefit, but at lower initial solids a static mixer was not needed. Finally, taking all the variables in consideration, it was found that a PEO dosage of 0.02-0.14 lb/1,000 gal was required to produce a dewatered product of 33 to 43 pct solids and screen underflow with turbidity of 20 to 50 NTU's, Table 1.

Table 1.--Results of the Alaska field dewatering test on placer effluents using 0.01 pct PEO

Mixer length, ft	Feed turbidity, NTU	Initial solids, pct	PEO dosage, lb/1,000 gal	Solids content, pct	Underflow, turbidity, NTU
600*	26,500	4.41	0.14	42.9	46
600	26,000	5.40	.19	40.5	45
600	14,000	2.42	.07	33.8	47
600	5,800	2.42	.06	42.2	33
600	5,400	.78	.05	42.5	31
600	5,200	.60	.03	20.1	40
600	2,000	.25	.02	35.6	32
600*	1,300	.52	.04	33.2	20
600	300	.09	.02	33.2	26
800*	15,000	2.68	.12	41.3	50
800	6,200	.68	.06	39.7	39

*Two element static mixer was used with pipe.

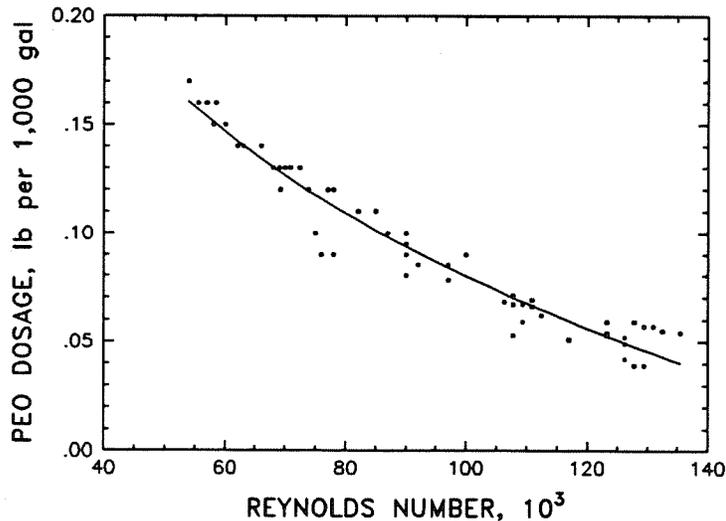


FIGURE 4.--Effect of Reynolds number on PEO dosage for the large field test unit.

A cost estimate of the Bureau of Mines process for dewatering Alaskan placer effluent streams with PEO was determined and is reported in an open file report (6). The cost estimates are for dewatering plants processing 1,000 gpm of effluent slurry at three representative turbidity levels. Both the placer and dewatering plants operated on the same 1 shift-per-day, 6 days-per-week schedule for the 100-day Alaskan operating season. Estimated fixed

capital costs for these plants processing placer slurries with effluent turbidities of 1,000, 3,000, and 5,000 NTU are approximately \$29,000, \$31,000, and \$34,000, respectively, on a fourth quarter 1986 basis (6). Operating costs are estimated to be \$0.34, \$0.37, and \$0.40 per thousand gallons of effluent slurry including amortization, plan and chemical cost.

The field testing at the crushed stone operation was conducted using the mobile unit previously described. As part of the operation for the stone quarry, a sand screw is used to remove coarse particles from the waste slurry prior to impoundment. The feed for the unit was taken from the pipe going to the impoundment from the sand screw. A wide variety of parameters was investigated. Polymer concentration ranged from 0.02 to 0.10 pct, Reynolds number from 4,800 to 16,000, and retention time from 30 to 136 sec. It was determined that a polymer concentration of 0.04 pct gave the best results. It was also determined that the flow of solids down the screen using the polyacrylamide was quite different from the previous tests where PEO was used. The PEO flocs flowed easily off the screen whereas, the polyacrylamide flocs tended to abrade as they moved down the screen causing some blinding of the screen. The difference between the two polymers is that the PEO has inherent antifriction properties allowing the flocs to move down the screen without abrasion. To overcome this problem, it was determined that tapping the screen periodically reduced the blinding allowing the screen to function. In practice this can be accomplished with an automatic tapping device set on a timer. The major parameter that had to be optimized in the crushed stone slurry was the mixing requirement. For this system it was determined that a Re of 4,800 with a mixing retention time of 68 sec produced 50 pct solids at a polymer dosage of 0.30 lb/1,000 gal treated. The material continued to dewater after exiting the screen and reached solids contents as high as 70 pct in 24 to 48 h.

At the present time, the polymer used in the study on the slurry from the crushed stone operation costs approximately \$0.50/lb when bought in bulk quantities. Therefore, the cost estimates for a crushed stone operation were estimated to be \$0.15 to \$0.25 per thousand gallons. The estimates are based on the operation where both a low cost anionic copolymer of acrylamide and a medium cost anionic polyacrylamide were used. Both polymers produced a dewatered product of 45-50 pct solids. The dewatered material has potential for use in a number of products, and any use would be a credit, this being the overall cost.

CONCLUSIONS

The field testing of two different waste slurries, one from a placer operation and the other from a crushed stone operation, has shown that the Bureau-developed polymer/static screen dewatering technique can be successfully used at low cost. To optimize the technique, the Re number and time of mixing in the feed pipe must be studied to obtain the optimum conditions.

REFERENCES

1. Scheiner, B. J., A. G. Smelley, and D. A. Stanley. Dewatering of Mineral Waste Using the Flocculant Polyethylene Oxide. BuMines B 681, 1985, 18 pp.
2. Sharma, S. K., and B. J. Scheiner. Effect of Physics - Chemical Parameters on Dewatering: A Case Study. Fluid/Particle Separation Journal, vol. 4, No. 3, September 1991, pp. 162-166.
3. Barksdale, R. D., Editor. The Aggregate Handbook. National Stone Association, Washington, DC, 1991.
4. Smelley, A. G., and B. J. Scheiner. Synergism in Polyethylene Oxide Dewatering of Phosphatic Clay Waste. BuMines RI 8436, 1980, 18 pp.
5. Scheiner, B. J., and M. M. Ragin. Factors Affecting the Dewatering of Phosphatic Clay Waste Slurries. Trans. of Soc. Min. Engr., vol. 284, 1988, pp. 1801-1805.
6. Magyar, M. J. Cost Estimate for Dewatering Alaskan Placer Effluents with PEO. BuMines Open File Report 39-87, 1987, 14 pp.

6/11/17

**INFORMATION AND COMMUNICATIONS PART 1:
HIGH-PERFORMANCE COMPUTING
AND NETWORKING**



USE OF HIGH PERFORMANCE NETWORKS AND SUPERCOMPUTERS FOR REAL-TIME FLIGHT SIMULATION

Jeff I. Cleveland II
Project Engineer

National Aeronautics and Space Administration
Langley Research Center
Hampton, Virginia 23681-0001

513-09
150483
p. 10

ABSTRACT

In order to meet the stringent time-critical requirements for real-time man-in-the-loop flight simulation, computer processing operations must be consistent in processing time and be completed in as short a time as possible. These operations include simulation mathematical model computation and data input/output to the simulators. In 1986, in response to increased demands for flight simulation performance, NASA's Langley Research Center (LaRC), working with the contractor, developed extensions to the Computer Automated Measurement and Control (CAMAC) technology which resulted in a factor of ten increase in the effective bandwidth and reduced latency of modules necessary for simulator communication. This technology extension is being used by more than 80 leading technological developers in the United States, Canada, and Europe. Included among the commercial applications are nuclear process control, power grid analysis, process monitoring, real-time simulation, and radar data acquisition. Personnel at LaRC are completing the development of the use of supercomputers for mathematical model computation to support real-time flight simulation. This includes the development of a real-time operating system and development of specialized software and hardware for the simulator network. This paper describes the data acquisition technology and the development of supercomputing for flight simulation.

INTRODUCTION

NASA's Langley Research Center (LaRC) has used real-time flight simulation to support aerodynamic, space, and hardware research for over forty years. In the mid-1960s LaRC pioneered the first practical, real-time, digital, flight simulation system with Control Data Corporation (CDC) 6600 computers. In 1976, the 6600 computers were replaced with CDC CYBER 175 computers. In 1987, the analog-based simulation input/output system was replaced with a high performance, fiber-optic-based, digital network. In 1990, action was begun to replace the simulation computers with supercomputers.

The digital data distribution and signal conversion system, referred to as the Advanced Real-Time Simulation System (ARTSS) is a state-of-the-art, high-speed, fiber-optic-based, ring network system. This system, using the Computer Automated Measurement and Control (CAMAC) technology, replaced two twenty year old analog-based systems. The ARTSS is described in detail in references [1] through [6].

An unpublished survey of flight simulation users at LaRC conducted in 1987 projected that computing power requirements would increase by a factor of eight over the coming five-years (Figure 1). Although general growth was indicated, the pacing discipline was the design testing of high performance fighter aircraft. Factors influencing growth included: 1) active control of increased flexibility, 2) less static stability requiring more complex automatic attitude control and augmentation, 3) more complex avionics, 4) more sophisticated weapons systems, and 5) multiple aircraft interaction, the so called "n on m" problems.

Having decided to continue using large-scale general-purpose digital computers, LaRC issued a Request for Proposals in May, 1989 and subsequently awarded a contract to Convex Computer Corporation in December of that year. As a result of this action, two Convex supercomputers are used to support flight simulation. The resulting computational facility provided by this contract is the Flight Simulation Computing System (FSCS). This system is described in references [8] through [11].

ADVANCED REAL-TIME SIMULATION SYSTEM

Through design efforts by both LaRC design engineers and design engineers at KineticSystems Corporation, three components of the ARTSS were developed to meet LaRC requirements. These were the serial highway network, the network configuration switch, and the signal conversion equipment. A block diagram of the ARTSS is presented in Figure 2.

Serial Highway Network

The LaRC ARTSS employs high-speed digital ring networks called CAMAC highways. At any given time, four totally independent simulations can be accommodated simultaneously. The equations of motion for an aircraft are solved on one of the mainframe computers and the simulation is normally assigned one highway. The purpose of the network is to transfer data between the central computers and simulation sites (control console, cockpit, display generator, etc.). The elements of a CAMAC highway are: the Block Transfer Serial Highway Driver (BTSHD); the fiber-optic U-port adaptor, the Block Transfer Serial Crate Controller (BTSCC); the List Sequencer Module (LSM); and the CAMAC crate. Three features of the networks were developed to meet the LaRC requirement. First, the mainframe computer interface to the BTSHD was developed. Second, the block transfer capability was developed to meet LaRC performance requirements. This capability resides in the BTSHD, BTSCC, and LSM. Third, the fiber optic capability was developed to satisfy our site distance problem. The simulator sites are from 350 to 6,000 feet from the computer center.

Prior to the development of the block transfer capability, a CAMAC message was approximately 19 bytes long which included addressing, 24 bits of data, parity information, and response information. The addition of the block transfer capability allowed for the inclusion of many CAMAC data words in a single message. During block transfers, data reads or writes proceed synchronously at one 24-bit CAMAC data word per microsecond. This is several times faster than the normal single word message rate. Besides the CAMAC standard message, there are two modes of block transfer. In the first, the entire block of data goes to a single module within a crate. It is implemented by the BTSCC repeating the module-select and function bits on the crate dataway for each CAMAC word. In the second block transfer mode, the block, either on read or write, is divided among several modules within a crate. This mode employs the LSM module which is loaded by the mainframe computer at set-up time with up to four lists of module-select and function bits. When this type of block transfer is in progress, the BTSCC acquires module number, function, and subaddress for each sequential CAMAC word in the block from the indicated list in the LSM.

To support Convex supercomputers, a new generation serial highway driver was developed. This driver provides direct connection to VMEbus and allows data to be streamed onto the network. Previous equipment transmitted 24 bits out of the 32 available on the host computer interface; however, the new hardware transmits the full 32 bits from the host computer. Packing/unpacking operations are no longer required to provide the 24 bits in 32 which results in lower input/output latency and increased computer time available for model computation.

Network Configuration Switch

The purpose of the network configuration switch system is to provide complete connectability between the simulation applications on the mainframe computer and the various simulation sites. Upon request, any sensible combination of available sites can be combined into a local CAMAC ring network in support of a single simulation. This network configuration for a given simulation is done during the initialization phase, after a highway has been assigned by the scheduling software. The application job requests sites by resource request statements and if the sites are available, the switch will electrically and logically configure the network without disturbing other running simulations. The switch is built for a maximum of 12 highways and 44 sites. Each highway may be connected to a different host computer. During the transition period, four computers were routinely used simultaneously doing flight simulation. Two of these computers were CDC CYBER 175 computers and two Convex supercomputers. In the final configuration, two Convex supercomputers with a total of six configuration switch ports are used.

Signal Conversion Equipment

Three types of output converter modules and two types of input modules were designed and built to LaRC specifications. The converters are high quality and have been added to the vendor's catalog. The digital-to-analog converters (DAC), analog-to-digital converters (ADC), and digital-to-synchro converters (DSC) are 16-bit devices with 14 bits of accuracy. The data transmitted uses 16 bits although only 14 are meaningful. This implementation allows LaRC to change converter precision without major changes in software or protocol. To decrease transmission time, data words are packed such that three converter words (16 bits each) are contained in two CAMAC words (24 bits each). The discrete input converters contain 48 bits per module and the discrete output modules contain 24 bits per module.

Clocking System

Flight simulation at LaRC is implemented as a sampled data system. The equations of motion are solved on a frame-by-frame basis using a fixed time interval. To provide the frame interval timing signals and a clocking system for synchronization of independent programs, LaRC designed and built the real-time clock system. This system is patented and is described in reference [7]. The clock system is composed of a central unit and multiple CAMAC modules called Site Clock Interface Units (SCIUs) which are connected by means of a separate fiber optic star network. Two distinct time intervals are broadcast by the central unit on a single fiber. The first time interval has a constant 125-microsecond period. The tic count necessary for a real-time frame is set in the SCIU by initialization software. This count is decremented by one for each occurrence of the interval timer. When the count reaches zero, each SCIU issues a signal that indicates beginning of frame. The frame time is determined independently for each simulation but must be a multiple of 125 microseconds. The second clock signal, called the job sync tic, has a longer period called the clock common multiple which is set manually, typically in the 1 to 3 second range. This longer period is used for synchronization. Each frame time must divide evenly into the clock common multiple, ensuring that all simulations will be synchronized on the occurrence of the job sync tic.

Figure 3 shows some of the wide variety of modules that can be accommodated in a CAMAC crate. The modules that are indicated as being developed to NASA specifications are available as standard catalog products from Kinetic Systems Corporation. The modules indicated as NASA developed were developed at NASA Langley Research Center.

REQUIREMENTS FOR NEW SIMULATION COMPUTING SYSTEM

The results of the 1987 survey of simulation users and program managers led to the definition of requirements for replacing existing computers.

CPU Performance

Real-time flight simulation at LaRC requires high scalar CPU performance to solve the equations of motion of the system being simulated. Using an existing simulation of an X-29 aircraft as a benchmark, the following CPU performance was specified:

1. If a single CPU configuration is provided, the CPU must solve the benchmark in at most 165 seconds.
2. If a multiple CPU configuration is provided, each CPU must solve the benchmark in at most 330 seconds.

Due to secure processing requirements, a minimum of two and a maximum of four independent computers were required. The CYBER 175 computer solves the benchmark in 660 seconds. Thus, the capabilities of the resultant total system will provide at least eight times the CPU processing power of the coupled CYBER 175 computers.

Real-Time Input/Output

The ARTSS CAMAC system has provided LaRC with a high performance real-time input/output system that has

extended the capabilities of the LaRC simulation system. Since ARTSS provides a high transfer rate with low latency, LaRC required provision of a compatible interface between the simulation computing system and the ARTSS CAMAC system. LaRC required that the new system include all software and hardware to connect to the ARTSS CAMAC real-time network. This connection was required to transfer block data over the network at a sustained rate of 24 million bits per second in the enhanced serial mode.

Responsiveness

One of the critical requirements for any real-time simulation system is system responsiveness. The FSCS system is required to respond to an external event, cause a short FORTRAN program to execute, and post an observable output response, in less than 150 microseconds. This elapsed time, called time-critical system response, is measured at an external port on the computer. The external event occurs at a repetitive rate of 1000 events per second. In addition to the time-critical system response, CAMAC input/output response is required to be less than 200 microseconds. CAMAC input/output response is defined as the time between the action of an interrupt generated in a CAMAC crate, transfer of one CAMAC word of data, execution of the short FORTRAN program, and transfer of one CAMAC word of output.

Frame Rate

To support simulation applications needing higher frame rates, LaRC required the system to run simulations at 1000 frames per second. At this frame rate, during any given frame, the system must deliver at least 600 microseconds of CPU time for the simulation model with 100 bytes of real-time input and 100 bytes of real-time output. The sum of system overhead and real-time input/output must be less than 400 microseconds.

Real-Time Data Recording and Retrieval

To support real-time data recording and retrieval during synchronous flight simulation, LaRC required the capability to record and/or retrieve information from two files for each simulation. The aggregate storage capacity was required to be a minimum of 180 megabytes. Sufficient data rate was required to permit a simulation to record or retrieve one 1000-byte record per real-time frame from each file simultaneously at a frame rate of 100 frames per second.

Language and other factors

At LaRC, almost all simulation programs are written in the FORTRAN language. Furthermore, simulations have been developed on CDC 6000 series computers and succeeding generations for over twenty years. With simulations written taking advantage of the CDC 60-bit architecture, LaRC required that the FSCS system support simulations written in the FORTRAN language with a minimum floating point mantissa precision of 14 decimal digits and with a minimum exponent range of plus and minus 250 decimal. In addition, the C language is required to support a limited number of applications, and Pascal is required to support the CAMAC configuration database.

An application development capability was required to operate simultaneously with simulations operating in real-time using all the real-time computing power specified. This application development capability has a minimum performance specification and required an advanced source language level debugger.

NEW SIMULATION COMPUTING SYSTEM

The computers that LaRC has put in place to fulfill the requirements are Convex Computer Corporation C3200 and C3800 series computers. These computers are classified as supercomputers and support both 64- and 32-bit scalar, vector, and parallel processing. The first delivery consisted of a Convex C3230 (3 CPUs expandable to 4) with two CAMAC interfaces. The system was delivered with two peripheral buses (PBUS): one PBUS that is used for input/output to standard peripherals such as tape, disk, and line printer and one PBUS that is used exclusively for real-time input/output to the ARTSS CAMAC network. Each VME Input/Output Processor

(VIOP) is a Motorola 68020 microcomputer that provides programmable input/output control. Each VIOP is connected to a standard 9U VMEbus and to the corresponding PBUS. The CAMAC interface consists of a KineticSystems Model 2140 Enhanced Serial Highway Driver for VMEbus. The second delivery consisted of one Convex C3850 (5 CPUs expandable to 8) computer configured similar to the C3230 with 3 PBUSs and two CAMAC interfaces. The computer contains 512 megabytes of main memory and sufficient disk and other peripherals to support flight simulation. The resulting computer configurations are shown in Figure 4.

There are four critical aspects of a computing system to support real-time simulation. These are: CPU performance, memory capacity, time-critical system response, and deterministic system performance.

The first computer installed (C3230) performs a simulation of an X-29 aircraft in 245 seconds per CPU which is 2.7 times faster than the computers being replaced. With two CPUs available for real-time, this results in over 5 times the CPU performance. The second computer (C3840) performs the X-29 in 117 seconds per CPU which is 5.6 times faster than the computers being replaced. With four CPUs available for real-time, this results in over 18 times the CPU performance.

Memory capacity is more than adequate to meet the requirements. The expanded memory capacity, compared with the old system, has allowed LaRC researchers to greatly increase the complexity of the simulations. The increase in memory capacity, coupled with the increase in CPU performance, has led to much higher fidelity simulations. Memory capacity is high enough to permit its use for real-time data storage and retrieval. If the data requirements of the real-time simulations exceed the memory capacity, a disk spooler program will be developed.

Time-critical system response is a measure of how fast the computing system can respond to real-time events from outside the computing system. Time-critical system response on both the computing systems has been measured at 31 microseconds, which exceeds the LaRC requirement.

Deterministic system performance is a measure of how consistently on a frame-by-frame basis the computing system calculates the simulation model without any loss in synchronization with real-time. To use a computing system for real-time simulation, the system must be able to solve the model in a very nearly fixed amount of time, no matter what the demands on the system are for other computing. The C3850 performs simulation models with less than one percent variation in model computing speed. Modifications to the C3230 software are being done to improve its model computing behavior.

Operating System

Convex Computer Corporation offers two real-time operating systems. The operating system currently in use at LaRC requires one CPU for all non-real-time activity: editing, program compilation, and other UNIX activities. The other CPUs may be dedicated to real-time simulation. At the request of a real-time program, the program is locked down in memory to prevent page faults and the CPU or CPUS are dedicated exclusively to the real-time program.

The second real-time operating system incorporates a specially developed real-time kernel that the entire operating system is built upon. With this version of the real-time operating system, the UNIX operating system portion will be pre-empted by real-time requests and the response to real-time interrupts will be deterministic and very short. This version supports, on a single CPU, all activities of a normal UNIX operating system and also simultaneously supports real-time applications. This operating system requires special hardware that is not available in LaRC computers.

CONCLUSION

NASA Langley Research Center is completing the development of a system to simulate in real-time increasingly complex and high performance modern aircraft. Utilizing centralized supercomputers coupled with a proven real-time network technology, scientists and engineers are performing advanced research using flight simulation.

Hardware and software developed and concepts used are applicable to a wide range of commercial applications that require time-critical computer processing including process control, power grid analysis, process monitoring, radar data acquisition, and real-time simulation of a wide variety of systems.

REFERENCES

1. Crawford, D. J. and Cleveland, J. I. II, "The New Langley Research Center Advanced Real-Time Simulation (ARTS) System," AIAA Paper 86-2680, October 1986.
2. Crawford, D. J. and Cleveland, J. I. II, "The Langley Advanced Real-Time Simulation (ARTS) System," AIAA Journal of Aircraft, Vol 25, No. 2, February 1988, pp. 170-177.
3. Crawford, D. J., Cleveland, J. I. II, and Staib, R. O., "The Langley Advanced Real-Time Simulation (ARTS) System Status Report," AIAA Paper 88-4595-CP, September 1988.
4. Cleveland, J. I. II, Sudik, S. J., and Crawford, D. J., "High Performance Processors for Real-Time Flight Simulation," AIAA Paper 90-3140-CP, September 1990.
5. Cleary, R. T., "Enhanced CAMAC Serial Highway System," presented at the IEEE Nuclear Science Symposium, San Francisco, California, October 23-25, 1985.
6. ANSI/IEEE Standards 583, 595, and 675, Institute of Electrical and Electronic Engineers, 1976.
7. Bennington, D. R., "Real-Time Simulation Clock," LAR-13615, NASA Tech Briefs, June 1987.
8. Cleveland, J. I. II, Sudik, S. J., and Grove, Randall D., "High Performance Computing System for Flight Simulation at NASA Langley," AIAA Paper 91-2971-CP, August 1991.
9. Cleveland, J. I. II, "Application of Technology Developed for Flight Simulation at NASA Langley," presented at the Technology 2001 Conference, San Francisco, California, December 3-5, 1991.
10. Cleveland, J. I. II, "Use of Convex Supercomputers for Flight Simulation at NASA Langley," presented at the Convex Worldwide User Group Conference, Richardson, Texas, May 17-22, 1992.
11. Cleveland, J. I. II, Sudik, S. J., and Grove, Randall D., "High Performance Flight Simulation at NASA Langley," AIAA Paper 92-4179-CP, August 1992.

Flight Simulation Research Requirements

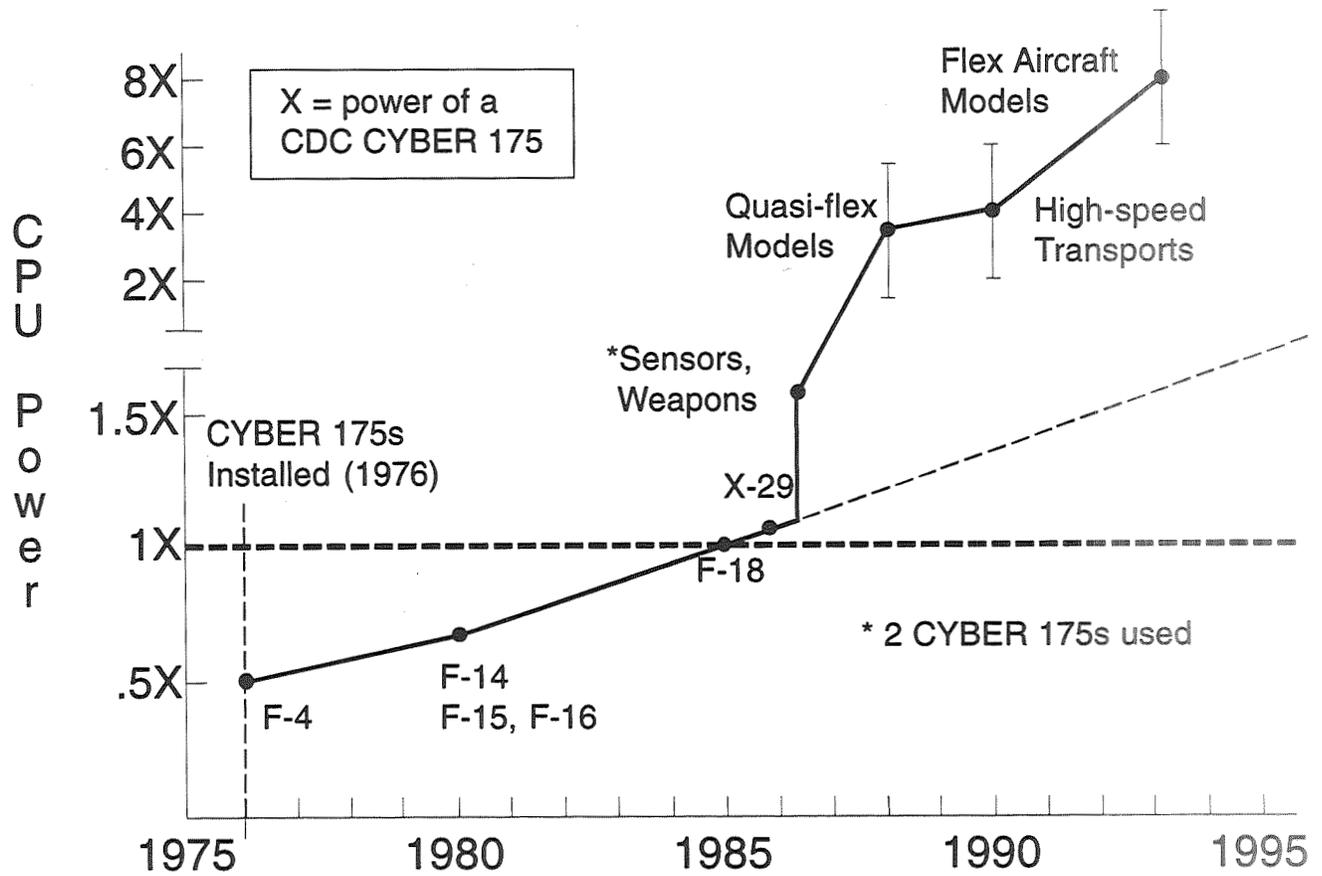


Figure 1.

Langley Flight Simulation System

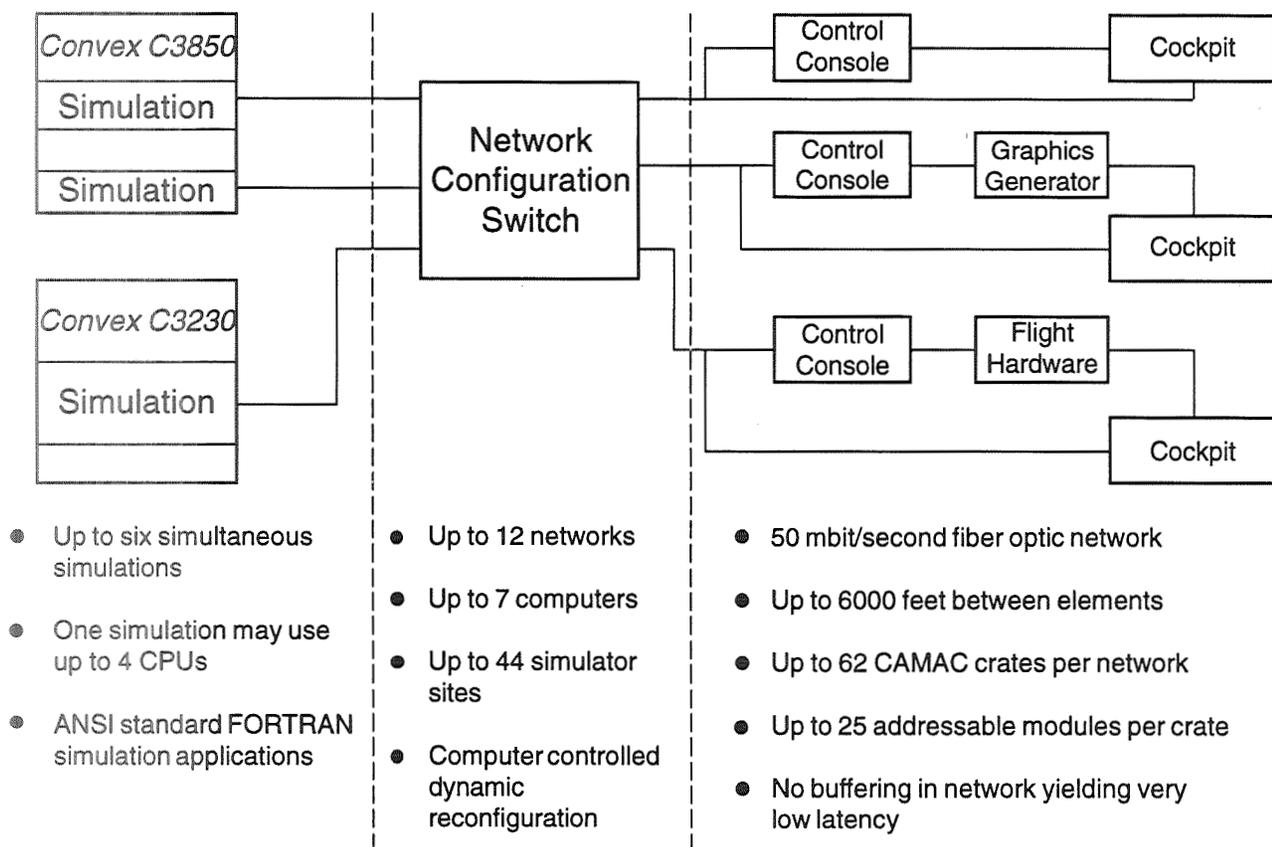
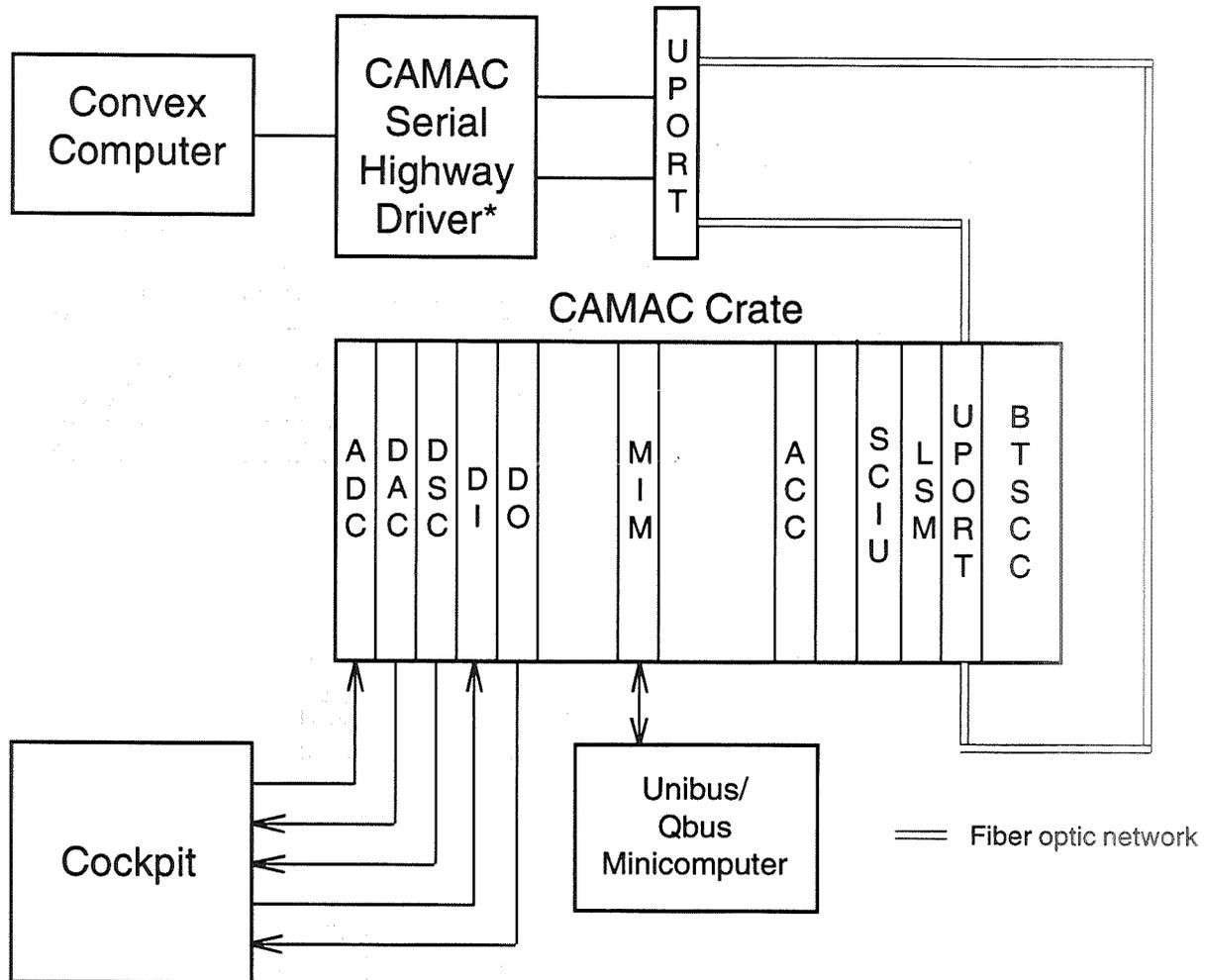


Figure 2.

Typical CAMAC Crate Configuration



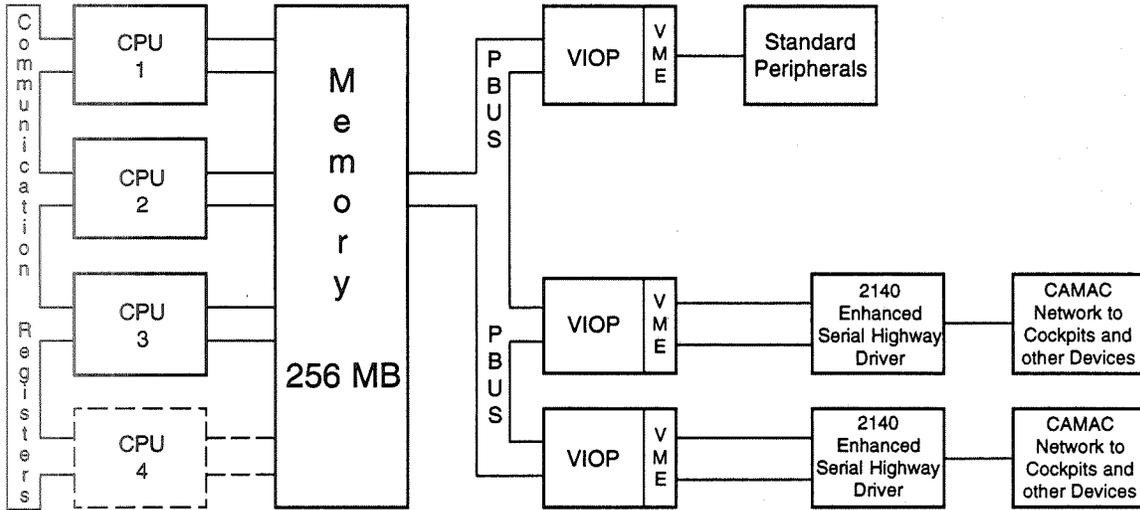
- ACC - auxiliary crate controller
- ADC - analog-to-digital converter *
- BTSCC - block transfer serial crate controller *
- DAC - digital-to-analog converter *
- DI - discrete input converter *
- DO - discrete output converter *
- DSC - digital-to-synchro converter *
- LSM - list sequencer module *
- MIM - minicomputer interface module **
- SCIU - site clock interface unit **
- U P O R T - fiber optic/electrical converter *

- * developed to NASA specifications
- ** developed by NASA

Figure 3.

Computing System Configuration

Convex C3230 Computing System



Convex C3850 Computing System

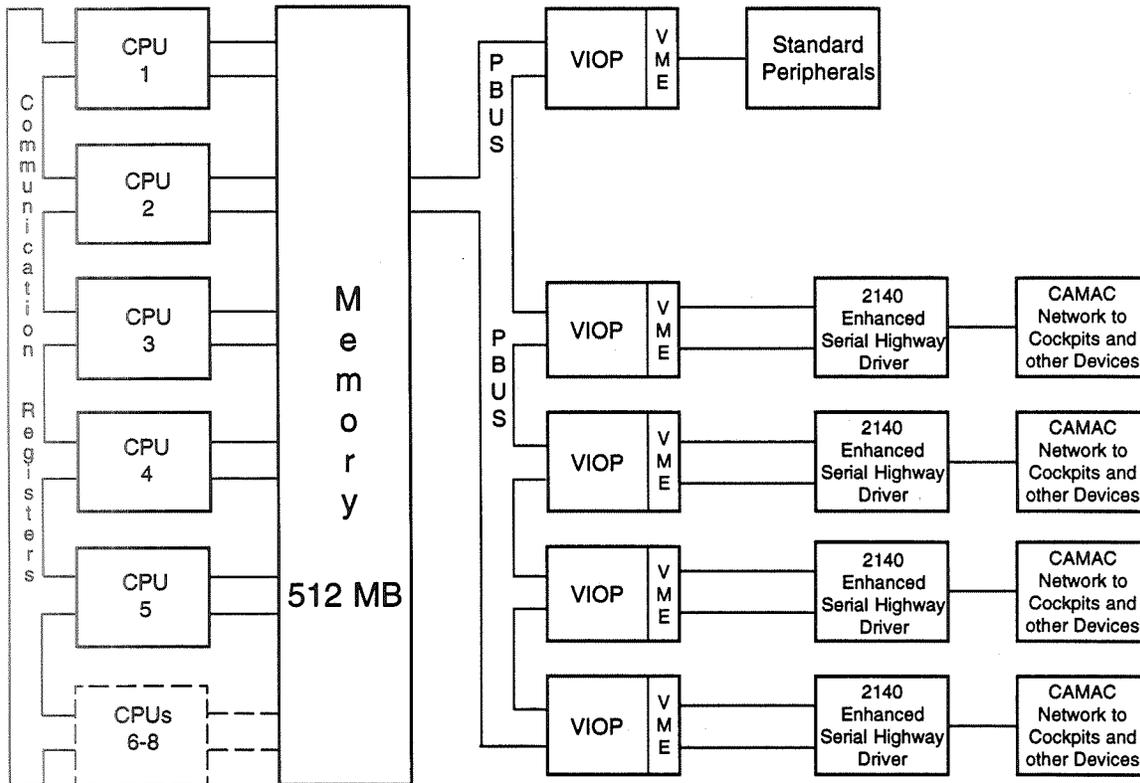


Figure 4.

SPLASH 2

Jeffrey M. Arnold
Duncan A. Buell
Walter J. Kleinfelder
Supercomputing Research Center
17100 Science Drive
Bowie, Maryland 20715-4300

314-61
150484
N93-25575

ABSTRACT

Splash 2 is an attached processor system for Sun SPARC 2 workstations that uses Xilinx 4010 Field Programmable Gate Arrays (FPGAs) as its processing elements. The purpose of this paper is to describe Splash 2. The predecessor system, Splash 1, was designed to be used as a systolic processing system. Although it was very successful in that mode, there were many other applications that were not systolic, but which were successful, nonetheless, on Splash 1, or that were not implemented successfully due to one or more architectural limitations, most notably I/O bandwidth and interprocessor communication. Although other uses to increase computational performance have been found for the Xilinx FPGAs that are Splash's processing elements, Splash is unique in its goal to be programmable in a general sense.

INTRODUCTION

Splash 2 is an attached processor system for Sun SPARC 2 workstations that uses Xilinx 4010 Field Programmable Gate Arrays (FPGAs) as its processing elements. The purpose of this paper is to describe Splash 2.

The predecessor system, Splash 1 [1], was designed to be used as a systolic processing system [2] [3]. Although it was very successful in that mode, there were many other applications that were not systolic, but which were successful, nonetheless, on Splash 1, or that were not implemented successfully due to one or more architectural limitations, most notably I/O bandwidth and interprocessor communication. Although other uses to increase computational performance have been found for the Xilinx FPGAs that are Splash's processing elements (see, for example, [4] or [5]), Splash is unique in its goal to be programmable in a general sense.

THE HARDWARE

The architecture of Splash 2 is shown in Figures 1 and 2.

The Splash 2 System

The system-level view of Splash 2 is shown in Figure 1. (This shows a 3-board system; a system can contain 1 to 13 boards.) An interface board plugs into the backplane and an SBus adapter board plugs into a Sun SPARC 2 workstation to run the Splash 2 system via the interface board.

This paper is taken from a paper published in the Proceedings of the 4th Annual ACM Symposium on Parallel Algorithms and Architectures and copyrighted by the Association for Computing Machinery. It is published here by permission of the ACM.

The interface board extends the address and data buses from the Sun into the address/data buses in the backplane. The Sun can read from and write to memory and memory-mapped control registers on the Splash 2 boards via these buses. The Sun provides only 25 address bits (that we take to be 23 since we deal only with data on 32-bit-word boundaries), which is inadequate to address the $13 \text{ (boards)} \times 17 \text{ (memories)} \times 512K \text{ (bytes)}$ of Splash 2 memory, so the interface board contains an address bank register that selects the Splash 2 board in the system.

There are three data paths into the Splash 2 system.

- (1) On the memory bus, data can be read and written into memories attached to each Xilinx processing chip.
- (2) A "linear data path" exists down the SIMD bus into the first Xilinx chip, X1, in the linear array of the first Splash 2 board in a daisy chain that can include as many as 13 Splash 2 boards. Output from the last Xilinx chip in the linear array, X16, of the first board passes as input to the X1 chip of the second board, and so on. Output from X16 of the last board in the daisy chain returns on the Rbus to the interface board.
- (3) A SIMD path exists by using the SIMD bus for broadcast. The SIMD bus has a data path into Xilinx chip X0 on each board, which can then inject SIMD instructions or data into the crossbar and thus broadcast to the other Xilinx chips on that board.

There are three modes for sending data into the Splash 2 system.

- (1) Splash 2 can communicate with the Sun via DMA transfers to and from the FIFOs of Figure 1. The two input FIFOs are $1K \times 36$ -bits; the two output FIFOs are $1K \times 32$ -bits. For these transfers, the interface board becomes a master on the Sun SBus and transfers bursts of data to or from the FIFOs. In typical operation the Sun programs and initializes the Splash 2 boards via memory-mapped transfers and then enables DMA for data transfer to/from Splash 2. In this mode, the 32 bits of data form the low 32 of the 36 bits in the FIFO. The high 4 bits are taken from a tag register. DMA data transfer can be sustained at about 38 Mbytes per second when the host workstation is CPU-loaded, or as fast as 54 Mbytes per second when the host is idle.
- (2) The Sun host can also perform direct writes to the input FIFOs of Splash 2. In this mode the high 4 bits of the 36-bit FIFO word are bits 5-2 of the address.
- (3) Splash 1 was and Splash 2 will be a useful processor for handling digital signals generated external to the Sun host. The external input accommodates input of such a signal directly to Splash 2. Further details of this are given below.

The Splash 2 Interface

The interface board is responsible for generating all the signals necessary in the backplane for running up to 13 Splash 2 boards.

The Sun data bus is latched and buffered to drive the backplane data bus for memory-mapped reads and writes. The Sun address lines are latched and buffered to feed the backplane, and the Sun can load an address bank register with a 7-bit address extension to obtain 30 bits of 32-bit-word addresses.

A clock generator provides the clock signal to the Splash 2 boards, can be programmed by the Sun to various frequencies, and can be programmed to single-step, N -step, or to stop on an interrupt.

Interrupts can be requested by any Splash 2 board, and the DMA controller can request an interrupt when transfers are completed. An interrupt register permits the Sun interrupt program to enable or disable interrupts and to read which interrupt source generated an interrupt. FIFO full/empty determination is under the control of Xilinx chips XL and XR.

The inclusion of Xilinx chips XL and XR was to provide for control of data transfer, clock (even a clock supplied by the external input), and tag bits independent of the Splash 2 boards. In Splash 1 such control was usually done in the first array chip, leading to asymmetry and crowded designs. With proper programming of XL and XR, the asynchronies of DMA transfer and external input and clock should not be seen by the Splash 2 boards themselves, and the XL and XR programs should function much like a system I/O library. A size register indicates the number of Splash 2 boards in a system, providing a signal to the Splash 2 boards so that one board is enabled to deliver data to the Rbus. A DMA controller performs SBus-compatible burst DMA transfers to and from the FIFOs in 16-word bursts.

To accommodate variable modes of data entry into a Splash 2 system, provision for an external signal input exists in the form of a daughterboard attached to the interface board. In this way, small changes in input signal conditioning can be made without requiring the entire board to be re-engineered. The daughter board can be configured to provide an external clock, thus allowing the Splash 2 system to be run synchronously with external data.

The Splash 2 Processing Boards

The Splash 2 board is detailed in Figure 2. Each board contains 17 Xilinx 4010 chips. Sixteen of these, X1-X16, form the processor array, connected both linearly and via the crossbar by 36-bit-wide data paths. The 17th chip, X0, has several uses to be mentioned later. Each of chips X1-X16 is connected via a 36-bit-wide path (18 address, 16 data, 2 control) to the $256K \times 16$ -bit memories. The memories can be read from or written to directly by the Sun on a 32-bit data path.

The linear data path brings data from either the previous Splash 2 board or from the SIMD bus into X1, through the linear array, and out from X16 to either the next Splash 2 board or to the Rbus, and from there to the interface board.

Among the many control lines on Splash 2 is a single interrupt line from each Xilinx chip back through the interrupt latch and mask to the host. This is useful for applications such as searches in which a Xilinx chip that found the solution can signal that fact back to the host and interrupt the processing. In addition, a global AND/OR and a global VALID line (GOR, GORV) extend from each Xilinx chip to the control chip X0, and a system global AND/OR runs from each Splash 2 board to the interface board.

A final feature of the Splash 2 board is the ability to load or store a configuration state into the Xilinx chips. Readout of the state was possible in Splash 1 and was invaluable for debugging and program optimization; the new ability also to load the Xilinx chips with a starting state configuration will greatly enhance the ability to monitor program behavior.

The 17th Xilinx chip X0 serves several functions. Its primary purpose is to control the crossbar. The crossbar itself is bit-sliced from nine TI SN74ACT884 4-bit crossbar chips. Up to eight different configurations can be chosen; X0 is used to select which configuration is in effect at any given cycle, and the crossbar control determines the direction in which data is transferred. Using multiple configurations can, for example, allow the 16 chips to be viewed as a two-dimensional mesh, or a 4-dimensional binary cube, provided that only one data path per Xilinx chip is used in any given cycle (since only one path exists). In this way, the realization of common communication patterns is relatively straightforward. For example, a 4-dimensional binary cube is realized as follows: View the linear array as a hamiltonian path through the 4-cube. Properly chosen, and with an appropriate coordinate labelling, this path provides all the connections in the x -dimension, four of eight in the y -dimension, and two and one, respectively, in the z - and w -dimensions. Six crossbar configurations, one for each direction for each of the y -, z -, and w -dimensions, now provide the additional connections to realize a 4-cube. Although arbitrary communication is not possible—only three and not four ports exist per chip—it is possible to communicate one dimension at a time, and many cube algorithms exhibit this characteristic pattern. In an analogous way one can realize a 4×4 mesh, although in one of the two dimensions only half the needed communication paths are available at a time.

A second function of X0 is to provide a broadcast capability into the crossbar. Splash 2 can be used as a SIMD computing engine, as will be discussed below, and the connection from the linear data path through X0 into the crossbar allows for a broadcast of instruction and immediate data to all chips on a board at a time, using the lines into the crossbar shared by X0 and X16.

To allow X0 to be sent "subroutine calls" in SIMD mode and to execute stored subroutines, and to allow for the lookup tables that can be expected to be heavily used, X0 possesses its own local memory.

One complication exists in that the memories are 16 and not 32 bits wide. To allow for both the host and the Splash 2 board to view the normal data width as 32 bits, the memories on the Splash 2 board are double-cycled; the host and the interface board pass 32 bit data to/from the Splash 2 board, and the board reads/writes 32 bits on word boundaries by using two cycles for every data transfer to/from the interface board. This design decision was based on the I/O pin count of the Xilinx 4010 chips. Many designs were considered, but it proved impossible to retain the linear array data path (2×36 bits), add a crossbar connection (1×36 bits), add a direct connection to memory (18 bits address, 32 data), and have any of the 160 I/O pins left over for control.

SIMD COMPUTING MODE

A Splash 2 board allows a 256-bit load/store in parallel to 16 Xilinx processing chips. The combination of crossbar and the linear array provides a powerful parallel data transfer capability similar to a network. With this view of the Splash 2 board, its use for SIMD computing is quite natural. To effect this mode of computing it is necessary to support broadcast of instructions and/or immediate data. This is made possible by lines running down the SIMD bus into Xilinx chip X0 of every board and from there directly to Xilinx chips X1 through X16 over the crossbar. In this mode, chip X0 could be programmed explicitly to serve as an instruction decode module and possibly also to convert from a vertical to a horizontal encoding of the instructions.

PROGRAMMING

There are three levels at which the Splash 2 system must be programmed: the Splash board, the interface, and the host. At the Splash board level the programmable components consist of the Xilinx processing chips, X1 through X16; the control chip, X0; and the crossbar. At the interface level the Xilinx chips XL and XR are user programmable. The host interface must provide input data streams and control the operation of the Splash system. A library of common control functions is provided for the interface board chips, XL and XR, and for the Splash board control chip, X0. Many linear and SIMD applications use only a single crossbar configuration, which can also be provided in a library. The host interface can be driven from either a C program that makes calls to a package of control routines or through an interactive graphical debugger. Therefore, the minimal Splash 2 program consists of a single replicated Xilinx program for X1 through X16 and a selection of library components for the rest of the system.

The programming environment for Splash 2 is based upon the VHSIC Hardware Description Language (VHDL). VHDL is a hardware specification language with many modern programming language features such as block structured control; user defined data types; and overloaded procedures, functions, and operators. VHDL programs can freely mix behavioral specifications with more traditional structural descriptions. The VHDL programming model includes the concept of time, so VHDL specifications can be simulated directly.

The Splash 2 programming methodology relies heavily upon simulation and logic synthesis. Users develop applications by writing VHDL behavioral models of their algorithms, which are then simulated and debugged within the Splash 2 simulator. Once an algorithm is determined to be functionally correct, it is compiled into a set of Xilinx chip configurations and the timing analyzed and optimized.

The Splash 2 simulator is a hierarchical model of the Splash 2 system comprising a set of VHDL models for each of the components of the system. When an application program is simulated, it is able to interact with the system exactly as it would with the physical hardware. The system models also verify that the application program meets any

hardware constraints such as memory sequencing and setup and hold times. Because the simulator is based upon commercial tools, a full source level debugging interface is available to the user.

A mix of logic synthesis and standard compilation techniques are used to compile VHDL programs into Xilinx configurations. A commercial logic synthesis tool is used to map the VHDL code into a gate list, where a peephole optimizer is used to perform a variety of Xilinx- and Splash-specific optimizations. The resulting gate list is then mapped into the CLBs and placed and routed using the Xilinx tool package. The Xilinx tools are used also to extract the detailed timing information from the placed and routed design. This information is used to construct a new VHDL model for each chip, which is then fed back to the Splash 2 simulator for timing analysis.

ACKNOWLEDGEMENTS

We acknowledge those who have contributed to Splash 2, including at least Neil Coletti, Steve Cuccaro, Elaine Keith, Brad Fross, Maya Gokhale, William Gromen, William Holmes, Daniel Kopetzky, Andrew Kopser, James Kuehn, Sara Lucas, Ronald Minnich, Michael Mascagni, John McHenry, Fred More, Louis Podrazik, Daniel Pryor, Craig Reese, Judith Schlesinger, Nabeel Shirazi, David Smitley, Douglas Sweely, Mark Thistle, Chris Tscharner, Paul Schneck, and Ken Wallgren.

REFERENCES

- [1] Maya Gokhale, William Holmes, Andrew Kopser, Sara Lucas, Ronald Minnich, Douglas Sweely, and Daniel Lopresti, *Building and using a highly parallel programmable logic array*, IEEE Computer **24** (1991), 81-89.
- [2] H. T. Kung, *Why systolic architectures?*, IEEE Computer **15** (1982), 37-46.
- [3] H. T. Kung and C. E. Leiserson, *Systolic arrays for VLSI*, Introduction to VLSI Systems, by C. A. Mead and L. C. Conway, Addison-Wesley, Reading, Massachusetts, 1980, pp. 271-292.
- [4] Will B. Moore and Wayne Luk (eds.), *FPGAs*, Abingdon EE & CS Books, Abingdon, England, 1991.
- [5] M. Shand, P. Bertin, and J. Vuillemin, *Hardware speedups for long integer multiplication*, Proceedings, ACM Symposium on Parallel Algorithms and Architectures (1990), 138-145.

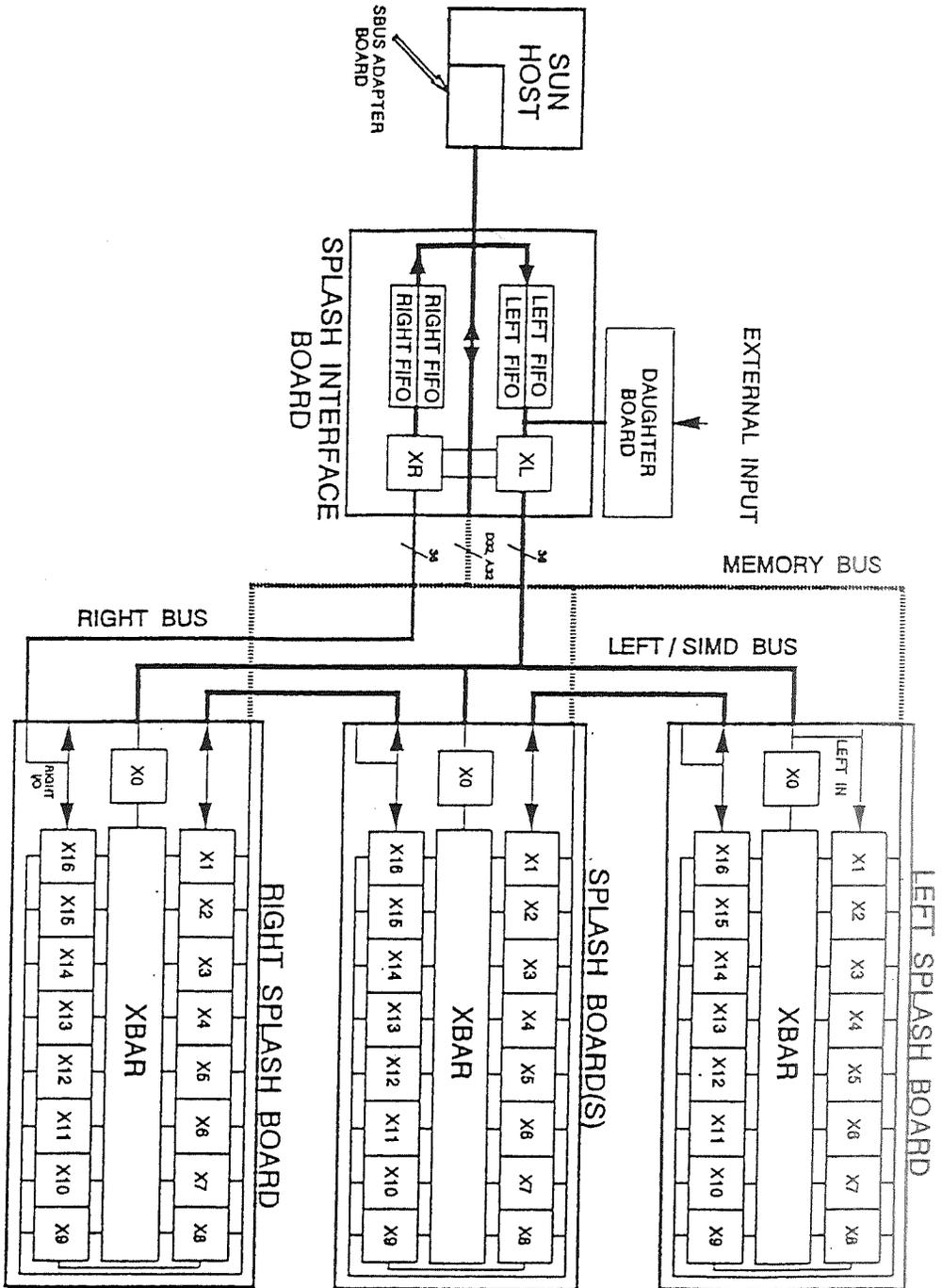


Figure 1: The SPLASH II System

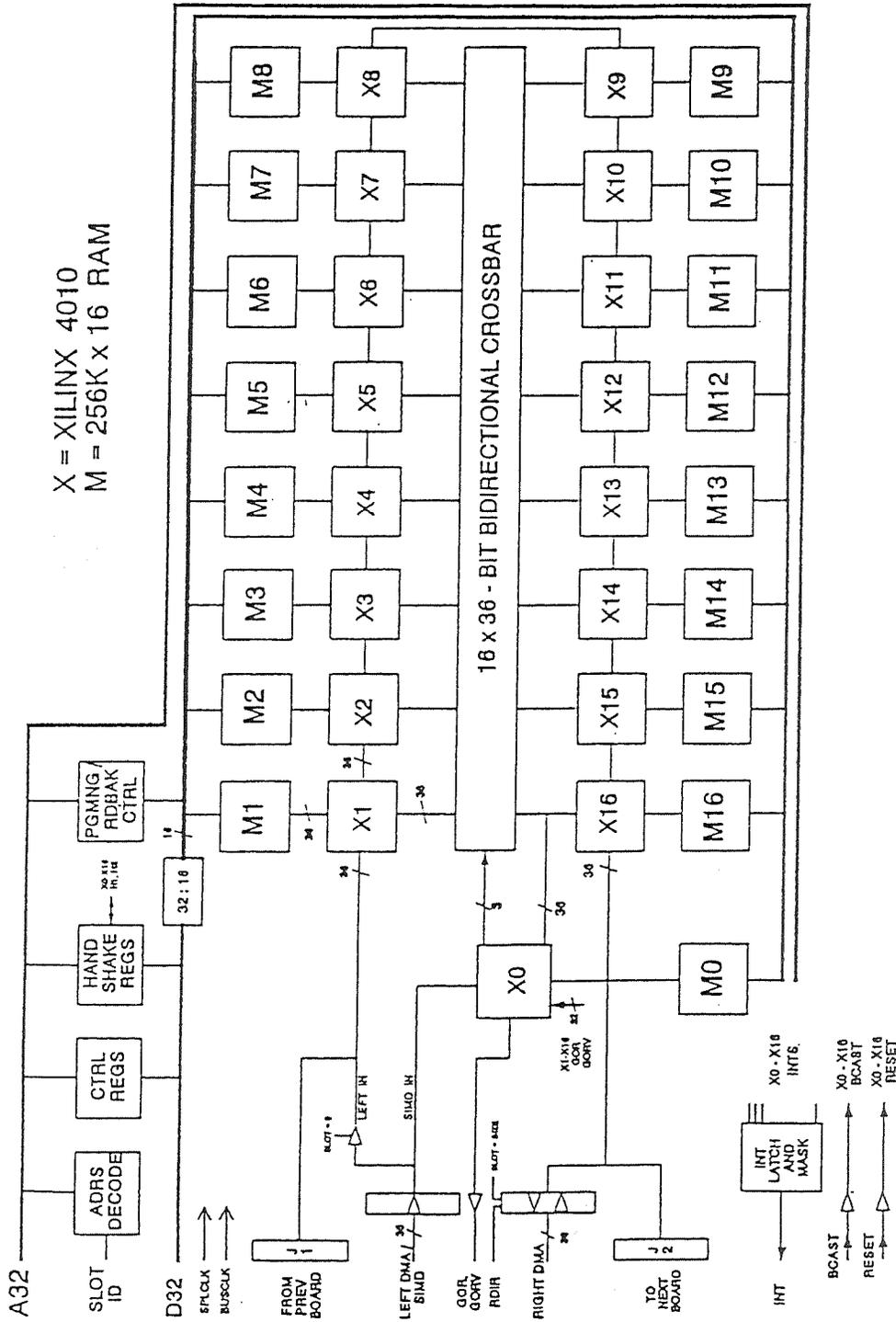


Figure 2: The SPLASH II Board

3/5-62

150485

P. 10

N93-25576

Object-Oriented Tools for Distributed Computing

Richard M. Adler
Symbiotics, Inc.
725 Concord Avenue
Cambridge, MA 02138

ABSTRACT

Distributed computing systems are proliferating, owing to the availability of powerful, affordable microcomputers and inexpensive communication networks. A critical problem in developing such systems is getting application programs to interact with one another across a computer network. Remote interprogram connectivity is particularly challenging across heterogeneous environments, where applications run on different kinds of computers and operating systems. NetWorks!TM is an innovative software product that provides an object-oriented messaging solution to these problems. This paper describes the design and functionality of NetWorks! and illustrates how it is being used to build complex distributed applications for NASA and in the commercial sector.

INTRODUCTION

A distributed computing system consists of software programs and data resources that are distributed across independent computers connected through a communication network [1]. Distributed systems are proliferating, owing to the availability of powerful, affordable microcomputers and inexpensive communication networks. Examples include: network operating systems for sharing file and printer resources; distributed databases; automated banking and reservation networks; office automation; groupware; and operations control systems for manufacturing, power, communication, and transportation networks.

In traditional uniprocessor systems, users connect with a mainframe or minicomputer from remote terminals to share its centralized data stores, applications, and processing resources. In distributed systems, both computing resources and users are physically dispersed across computer nodes on the network. Distribution offers important advantages, including replication, fault tolerance, and parallelism. For example, heavily-used programs or data may be duplicated on multiple nodes to increase resource availability and reliability. In addition, individual applications may be distributed, enabling their constituent tasks to be executed simultaneously on dedicated computers.

These benefits are achieved at the price of increased design complexity.¹ In particular, users and programs must frequently access information resources and applications that reside on one or more remote systems. Such access presupposes solutions to the design problems of:

- establishing interprogram communication between remote applications.
- coordinating the execution of independent applications, or pieces of a single distributed application, across networked computers.

Interprogram connectivity enables applications to exchange data and commands without outside intervention (e.g., users explicitly transferring and loading files). Such capabilities are critical for highly automated distributed systems that support activities such as concurrent engineering, data management and analysis, and process or workflow control. Connectivity also enables end-users to access remote database or file system resources interactively, via graphical user interface programs.

¹Obvious space limitations preclude discussion of important design issues such as: (a) locating distributed resources, which may be replicated and/or movable; (b) maintaining consistency across replicated data stores or distributed computations; and (c) managing recovery from partial failures such as node crashes or dropped network links in an orderly and predictable manner.

Coordination dictates how distributed applications interact with one another logically. Common models include client-server, peer-to-peer, and cooperative groups [2].

Three basic strategies have been developed to address the problems of distributed interprogram communication and control — network programming, remote procedure calls, and messaging systems. This paper focuses on NetWorks!TM, a novel messaging system that is based on advanced object-oriented software technologies. NetWorks! provides a network computing solution that is generic, modular, highly reusable, and uniform across different hardware and software platforms. This last attribute is important because many distributed computing tools are restricted to homogeneous environments, such as PCs or workstations. NetWorks! conceals the complexities of networking across incompatible operating systems and platform specific network interfaces, enabling application developers that lack systems programming expertise to build complex distributed systems.

The next three sections of this paper review and compare network programming, remote procedure calls, and conventional messaging systems. The fourth and fifth sections describe the NetWorks! system and illustrate its use through distributed applications in the domains of process control and office automation. The sixth section discusses extensions to NetWorks! that are currently being commercialized, which highlight the benefits of its innovative object-oriented approach.

NETWORK PROGRAMMING

Networked computers exchange data based on a common set of hardware and software interface standards, called protocols [3]. Examples include Ethernet, TCP/IP, LU6.2, and DECNet. The low level physical protocols consist of devices resident in computers connected together with a cable (e.g., Ethernet interface cards and thinnet cable). Programs called device drivers enable a computer's applications and operating system to communicate with its physical interface to the network. Higher level software protocols dictate how remote systems interact during the various phases of data exchange. For example, the receiving computer may need to acknowledge receipt of data to the sender. Both systems must agree on the convention (i.e., datum) used to signal acknowledgments.

Software network protocols have an advertised application programming interface (API), which takes the form of a library of function calls. Network programming consists of using protocol API libraries to define interactions between remote applications. Specifically, developers must insert API calls into programs to: initiate network connections between source and target programs across their respective host computers; send and receive the required data; and terminate connections. For example, developers might use network programming to connect diverse tools for computer-aided engineering. This would enable users to extract a design model from a database on one computer, apply a simulation program to the model on a second, analyze the results, and revise the design with another tool on a third node.

Network programming is complex and highly detailed: every phase of an interaction must be handled explicitly (e.g., managing network connections, the sequence of data exchanges that constitutes a complete "session," translating across the data and command interfaces of different applications). Network programming is particularly difficult across heterogeneous computers, requiring specialized expertise with incompatible protocols and operating systems. Typically, the code required to tie programs together is highly application-specific, which limits opportunities for reuse. Moreover, network programming code is generally tightly coupled to application-specific behaviors, impeding maintainability and extensibility of both the applications and their distributed computing interfaces.

REMOTE PROCEDURE CALLS

Remote Procedure Calls (RPCs) represent a second generation of tools for distributed computing [4]. RPCs are commonly used to implement client-server models for distributed interactions. One program, called a client, requests a service, which is supported by some other application, called a server. Upon receiving a

client request, the server responds by providing the requested service. Common examples of client-server computing based on RPCs include database, network file, and network directory services.

RPCs extend the familiar function call mechanism used on single computers to work over the network. In the single processor case, a program invokes a subroutine and blocks until the subroutine returns results and control. RPCs distribute this model by using special programs called "stubs" to link applications that reside on different computers. An RPC from an application is automatically directed to the relevant (local) client stub program. The client stub dispatches the call to the corresponding (remote) server stub, which invokes the target application. The latter eventually completes its processing and transfers results and control back to the server stub. This stub relays the results back to the client stub, which returns them to the original calling application, as depicted in Figure 1.

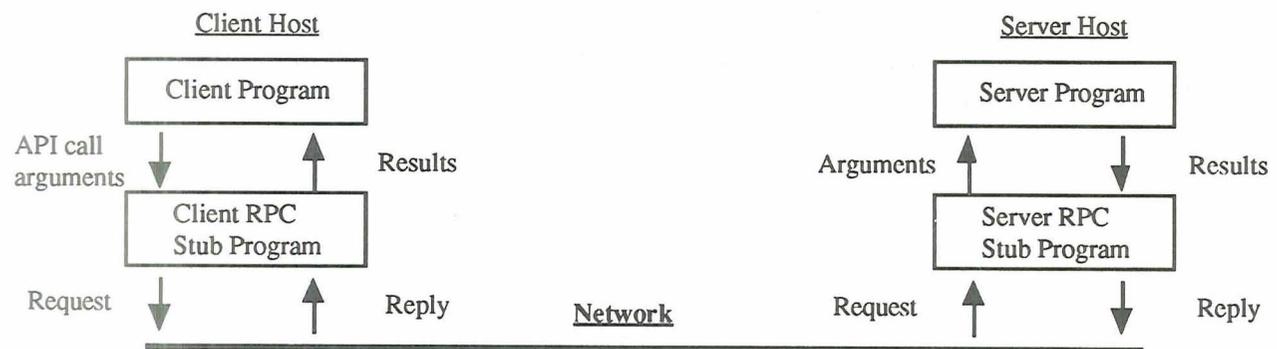


Figure 1. Client-Server Interaction using Remote Procedure Call Model

RPC tools generally include an API call library, a specification language, and a stub compiler. Developers insert the necessary RPC API calls into the source and target application programs. They then use the specification language to define the desired aspects of the interaction between the remote applications, such as packing and unpacking call parameter values and results, and managing network connections for transporting request and response data. Developers then compile their specifications to generate the RPC stub programs.

RPC tools represent a major improvement over programming with network protocols. First, application developers are already familiar with traditional, nondistributed procedure call models. Second, distributed behaviors are largely isolated within RPC stub programs, fostering modularity and maintainability. Third, stub compilers automatically generate code for interfacing with low-level network protocols. This greatly reduces the amount of programming required to define interactions between remote applications as well as the amount of systems level expertise required to design a distributed system. Fourth, RPCs specify program interfaces in terms of named operations with type signatures for call parameters. This means that RPC calls can be debugged for argument typing errors at compile time. Finally, RPC implementations are widely available and are often bundled into operating systems (e.g., Unix V4.2 BSD, the Open Software Foundation's Distributed Computing Environment).

RPC tools also have important limitations. First, RPCs are not fully compatible across disparate operating systems, network protocols, and tool vendors, although standards are emerging. Second, any distributed behaviors that cannot be expressed via the RPC specification language must be implemented using network programming. Intermingling custom code with the stub programs generated by the specification compiler complicates maintainability and extensibility. Third, the specifications for RPC stub programs are not readily reusable across applications. Fourth, every RPC-based service must be initialized explicitly and must always be executing, which wastes processing resources. Finally, the simple function call mechanism underlying RPCs is difficult to extend from client-server architectures to more advanced interaction models. In particular, standard RPCs are blocking or synchronous: such behavior is inefficient as it precludes calling programs from performing other tasks pending return of

results. RPCs are also inherently pairwise and unidirectional, making them unsuitable for distributed control models such as peer-to-peer, extended conversations, and cooperative groups.²

MESSAGE-PASSING MODELS

Message-passing models provide more flexibility for constructing distributed interactions than function call semantics. In particular, the act of passing a message need not commit the sending program to block pending replies. For example, a Request-Only model sends a single message from a source to a target process. It is used primarily for efficient, "side-effect" interactions, such as distributing information from data feeds. Similarly asynchronous (i.e., non-blocking) Request-Reply models support concurrent computing, in which programs dispatch service requests and immediately proceed with other activities pending the return of replies. The added reply message can be used to return either results of remote computations or acknowledgments that requested actions such as database updates have been completed. RPCs correspond to the synchronous Request-Reply model. Other messaging models include broadcasting (one-to-all), multicasting (one-to-many), and cooperative groups, which encompass specialized control protocols such as voting or negotiation. Note that the flexibility afforded by asynchronous models incurs distributed control overheads for tracking pending messages and responses.

Messaging systems come in two basic varieties, pipes and object-oriented. A pipe establishes a stream or channel between two applications [5]. Typically, pipe tools supply an API library for initiating pipe connections, writing to them, checking status, and reading from them. For example, data may be sent down a pipe from a source process and received at a sink process. All of the network level programming required to create pipes and transport messages through them is built into the pipes tool, and concealed from developers through the high level API. Pipes tools generally use the asynchronous Request-Only model, which is useful for efficient point-to-point transfer of bulk data.

Queue-based messaging represents a variation on pipes that establish connections between computers rather than specific programs. Such tools typically consist of two components, an API library and a queue management system. A client-server application would be implemented as follows (cf. Figure 2). First, the client program executes a Send API call, which posts a request message to the target server onto the local outbound queue. The local queue manager dispatches the message to the target node, where the remote queue manager posts it to an inbound queue. The server program uses a Receive API call to retrieve the request from the inbound queue, and then performs the required processing. It then uses a Send call, which posts the response to the outgoing queue. The remote queue manager sends this reply message back to the inbound queue on the client node, where the client program can retrieve it. Recently, directories have been incorporated into some messaging systems, which identify the nodes where registered services are located. This design feature enables client programs to issue service requests transparently, without having to know in advance or specify the server's whereabouts.

Message-based tools are much simpler to use than RPCs in that they dispense with stub specifications, compilers, and custom network programming. The messaging API calls completely conceal interfaces to network protocols, which fosters uniformity and portability of messaging tools across heterogeneous systems. However, conventional messaging systems fail to partition code for distributed computing in a manner analogous to RPC stubs. The logic for packing and unpacking message contents and for polling the pipe or queue to check message status is coded into the application program proximate to messaging API calls. Moreover, models for distributed interaction that are not directly supported (e.g., cooperative groups, multicasting, reliable transaction protocols) must be implemented using the available API library functions. A substantial amount of distributed computing code may need to be inserted into applications to realize the requisite logic, data management, and interprogram coordination. Although developers could apply structured design techniques, a more elegant strategy is to promote modularity

²A few RPC tools have been extended to try to address these objections (e.g., with callbacks or futures); however, such modifications violate the basic semantics of procedure calls.

within the framework of the distributed computing tool itself. Ideally, the tool would also foster reusability of distributed functionality, another feature that is lacking in conventional messaging tools.

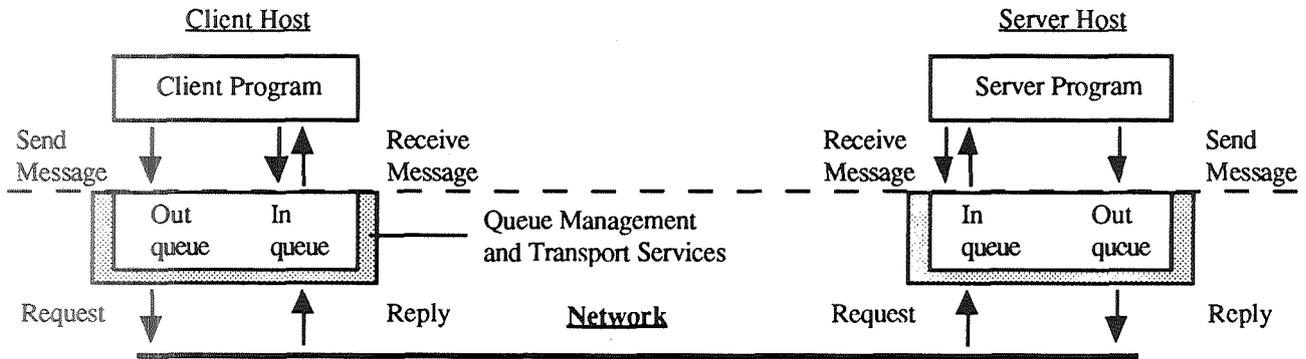


Figure 2. Client-Server Interaction using a Queue-Based Messaging Model

NETWORKS — OBJECT-ORIENTED MESSAGING

Like conventional messaging systems, NetWorks! defines a uniform high level communications interface between distributed programs running on heterogeneous platforms. In particular, a NetWorks! Messaging Facility (NMF) resides on all nodes of the network. The NMF provides queues, queue management, and message transport facilities. The NetWorks! architecture differs from conventional messaging systems in two primary respects. First, it introduces objects called Agents. Agents provide the locus for: (a) the control logic for passing and retrieving messages; (b) packing and unpacking application data for messages; and (c) distributed coordination. Intrusions into applications are thereby limited to high-level API calls, such as telling Agents to initiate messages for distributed interactions. Thus, instead of manipulating the NMF message queues directly, applications interact with Agents. In this respect, Agents play a role analogous to RPC stub programs. A basic NetWorks! model for a client-server interaction is depicted in Figure 3. More detailed examples are presented in the next section.

Second, the NMF incorporates a scheduler, which automatically manages incoming messages. Basically, the scheduler identifies the target Agent from each message and initiates the execution of that Agent for the given message. In essence, the scheduler delivers incoming messages to the relevant Agents, which then interact with their associated applications by injecting appropriate data or commands. This architecture greatly simplifies the control of asynchronous interactions.³ The NMF scheduler is basically nondeterministic, although the messaging API enables developers to specify values for a priority attribute to control precedence ordering of messages explicitly.

NetWorks! provides a development environment for defining Agent objects and for exercising and debugging their behavior interactively. An Agent consists of: (a) conventional program code, such as C or C++, (b) calls to the Agent API library for interacting with the NMF; and (c) calls to the native API of the Agent's associated application(s). Agents can be created on one platform and generated and installed (i.e., compiled and linked) remotely, across heterogeneous computers, for true distributed development. Developers then connect applications to Agents using a high-level messaging API, which is uniform across different programming languages. Agents process API calls from applications to prepare messages, which are posted to the NMF for transport. Agents also receive incoming messages and inject them directly into their associated applications, simplifying control for asynchronous interactions.

³For flexibility, the NetWorks! messaging API supports function calls to retrieve messages from Agents in situations where polling behavior is desired.

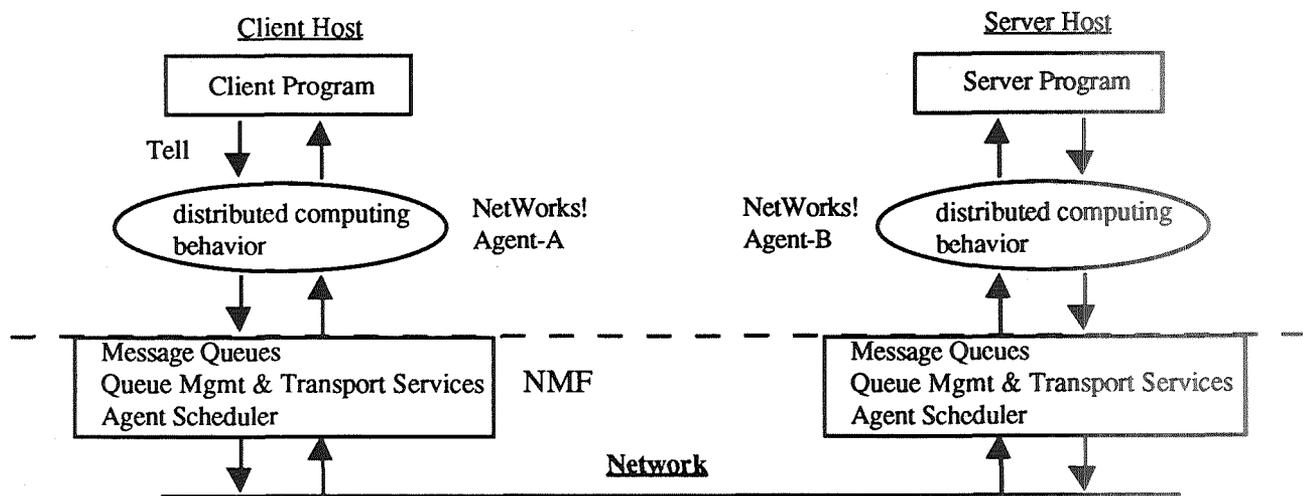


Figure 3. Client-Server Interaction using NetWorks! Object-Oriented Messaging Model

NETWORKS! APPLICATION EXAMPLES

Process Control

NetWorks! Agents support an object-oriented approach to remote interprogram connectivity. Objects are data structures that contain state information and behavior, which consists of operations for accessing and manipulating state information [6]. These operations, called methods, are invoked by sending messages to particular objects. This message-passing model maps directly into distributed interactions among application elements viewed as objects. Objects can be reused by creating new subclasses, which physically reuse or *inherit* behaviors from their parent classes. Moreover, inherited behaviors can be customized (i.e., overridden) selectively, through a process known as specialization. Object models encapsulate internal state and behavior by restricting any access other than through the message-based interface. This form of information-hiding or modularity enables system elements to be developed and tested independently, for later integration.

Consider a simple distributed system that uses NetWorks! Agents to support polling of remote sensors by an application that performs fault detection, isolation, and recovery functions (FDIR). FDIR functions are critical for automating operations support for mission control centers, manufacturing, or other complex network systems. The FDIR program initiates the polling sequence (cf. Figure 4) by issuing the NetWorks! application API call:

(Tell :agent FDIR-1 "poll measurement-Z")

The message contents in this case consists of the polling command for *measurement-Z*. The *:agent* keyword indicates that this message is to be delivered to the Agent *FDIR-1*, which is assumed by default to be co-resident with Program-A. The NetWorks! *Tell* API function is asynchronous, enabling the FDIR program to move on immediately to poll other sensors or perform FDIR reasoning. Agent *FDIR-1* may also support other messages from the FDIR program (e.g. for querying remote databases). The synchronous *Tell-and-Block* API function is provided to support blocking control models as well.

NetWorks! Agent objects contain two methods that control the processing of messages, called in-filters and out-filters. An in-filter parses incoming messages and then (a) invokes the Agent's associated application and/or (b) relays the message, which it may modify, to another Agent. Developers use the NetWorks! *Pass* API call to send messages from *within* an Agent in-filter method to another Agent.

(Pass :agent s-monitor :host host-2 "poll measurement-Z")

In this case, the *Pass* API function relays the polling request from Agent *FDIR-1* to the Agent *s-monitor*, which is located on the remote platform *host-2*. The *Pass* function posts the message to the local NMF, which transports it to the NMF on *host-2*. This NMF executes Agent *s-monitor*'s in-filter method, which parses the polling command message and interacts with the controller program for *sensor-Z* to collect the relevant measurement. Upon completing this process, the in-filter uses a *Set-Contents* API command to store the acquired data value in the *host-2* NMF queue and terminates. The *host-2* NMF then invokes the out-filter method for Agent *s-monitor*. Out-filters enable developers to postprocess results from in-filters. For example, a timestamp could be appended to the measurement value. The two NetWorks! NMFs automatically transport the final response message across the network to Agent *FDIR-1*, which injects it into the FDIR program using the latter's native API.

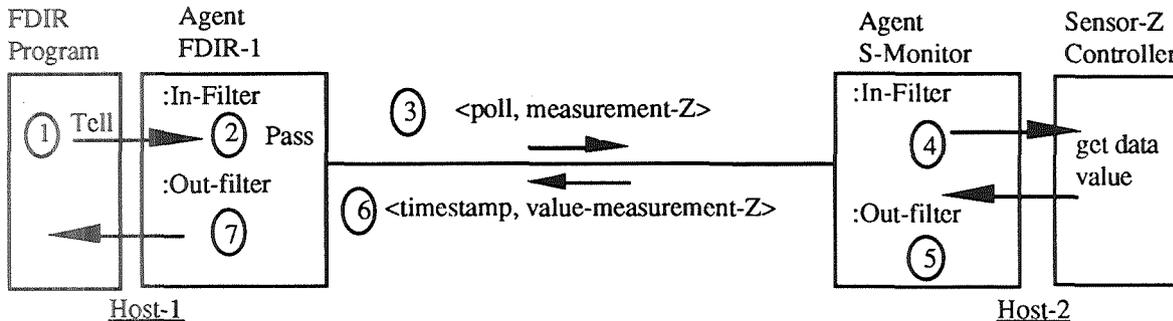


Figure 4. NetWorks! Agent-Based Message passing Model

NetWorks! Connect — DDE Over the Network

Connect is a commercial end-user software product that Symbiotics constructed using the NetWorks! development tool. Connect is a distributed application that extends Dynamic Data Exchange (DDE) over the network. DDE is a protocol within Microsoft's PC Windows operating environment that links data across two applications running on the same computer. These links are established using standard Windows application menu operations such as Copy (to the Clipboard) and Paste Link (from the Clipboard). The Paste Link operation links the copied data dynamically, so that updates made to the original source information (e.g., cells in a spreadsheet), are automatically propagated to update the linked information in some other application (e.g., a table in a word processing document).

NetWorks! Connect extends these operations to link data across DDE-compliant applications that run on different PCs across a network. Connect relies on NetWorks! Agent that copy remote Clipboards and handle DDE messages on each platform. The only departure from the single platform linking process is that users interact with a NetWorks! Windows-style menu to select a remote computer and retrieve the contents of its Clipboard. The DDE Agents transparently maintain the link produced by the Paste Link operation over the network. Connect also provides utility Agents for sending and receiving files across the network, and for exchanging *ad hoc* memos (e.g., for on line conferencing). In addition, "power" users can embed NetWorks! API calls into custom macros, Visual Basic programs, or programs in other languages that support DDE API calls. Lastly, because NetWorks! runs on heterogeneous platforms, such Connect applications can exchange data with programs on non-Windows platforms.

EXTENSIONS TO NETWORKS! — THE AGENT LIBRARY

An important advantage of NetWorks! over conventional distributed computing tools is that it promotes reusability through object-oriented inheritance of Agent behaviors. To exploit this advantage fully, Symbiotics is developing a library of Agent classes that establish predefined models for integration (Gateways) and distributed control (Managers). This library will facilitate a highly intuitive building block approach to designing and implementing distributed systems: developers can

simply select suitable library Agents and specialize them to satisfy application specific requirements. The Agent library will be commercialized next year as a layered NetWorks! product.

Gateway Agents

The Gateway Agent class defines a uniform peer-to-peer interaction model for integrating heterogeneous programs. Peer-to-peer simply means that any application Agent can act as a client or as a server with respect to any other. This control model, which is implemented via in-filter and out-filter methods, invokes a set of auxiliary Agent methods in a data-driven manner to process:

- API calls from the Gateway's application to initiate service requests (client behavior).
- incoming request messages from other application Gateway Agents (server behavior).
- responses to prior outgoing service requests from the Gateway (client behavior).

These messages all conform to a standard format, which enables the Gateway control model to determine which type of behavior is required.

An application is integrated into a distributed system by creating a new Gateway subclass and specializing two sets of Agent methods. One set, called translator methods, maps information across incompatible data and command interfaces for independent applications. NetWorks! incorporates a Data Management Subsystem (DMS), which supports a uniform "neutral exchange" format for all Gateway messages. DMS defines an object-based model for representing both primitive data types (e.g., character, integer, float) and composite types (e.g., database records, frames, arrays, files). A high-level API enables developers to create new composite types, and to create, pack, and unpack instances of these primitive and composite types. (Note: RPCs often offer similar tools, such as XDR [7].) A Gateway Agent's translator methods use the DMS API along with the application's API to map between the neutral exchange and native formats. The second set of Gateway methods defines the program's desired client and server behaviors, using the translator methods as a high level API to move data and commands into and out of the application non-intrusively. These methods typically consist of program Case statements, whose individual clauses dictate specific client or server behaviors. Gateways thereby separate communication from data management, and cleanly partition both kinds of functionality from the behaviors required for an application to participate in a distributed interaction.

Gateway Agents promote reusability of code in two ways. First, every application Gateway subclass inherits the uniform peer-to-peer control model from the root Gateway Agent class. Gateways thereby conceal and reuse common message passing *and* control behaviors for inter-Agent communication. Second, applications are often implemented using a development tool, such as 4GLs, CASE products or AI shells. In these situations, translator methods can be defined once, in the tool-specific Gateway subclass. Subclasses of that Gateway then inherit both messaging and information mapping functionality, leaving the developer to specify only the desired client and server behaviors for particular applications. The uniformity and modularity afforded by Gateways promote maintainability and extensibility, which are particularly important to support large, complex distributed systems that evolve over extended lifecycles in the field.

Manager Agents

Gateway Agents integrate distributed applications, facilitating basic peer-to-peer interactions. Manager Agents define reusable control models for coordinating more complex interactions. For example, a Hierarchical Distributed Control (HDC) model decouples application Agents from direct connections with one another [8]. Application Gateways send request messages to the HDC-Manager, which posts them to a task agenda (cf. Figure 5). The HDC-Manager contains a directory that identifies the application Gateway, its location, and its request interface for each service that is available in the distributed system. Using this directory, the HDC-Manager processes its agenda asynchronously, dispatching requests to the supporting Gateway for each pending task. Each such Gateway posts service

results back to the HDC-Manager, which relays them directly back to the relevant requesting Gateway Agents. Using this model, application Gateways need only know (a) how to interact with the HDC-Manager, and (b) the services that are available through it — no prior knowledge is needed about the identity, location or access interface of any other Agents. This design also promotes maintainability and extensibility for complex distributed systems: the HDC-Manager directory creates a layer of control abstraction that insulates application Gateways from modifications to each other and additions of new application Gateways and services.

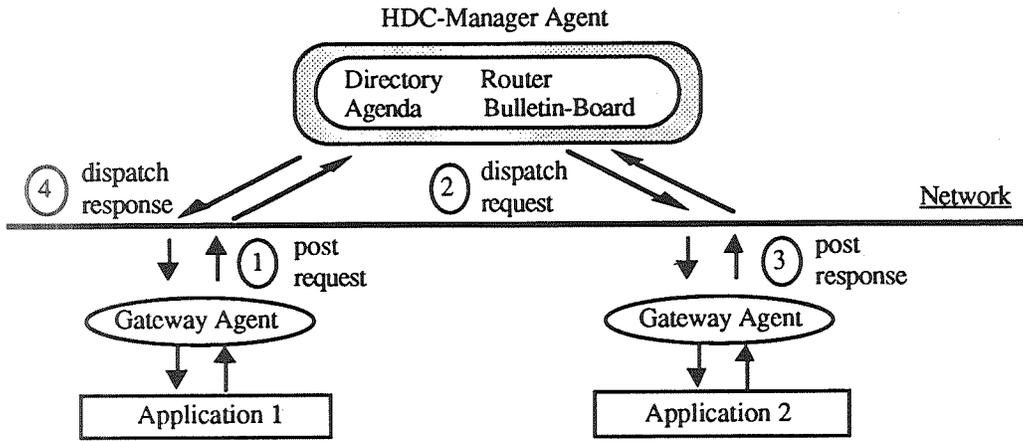


Figure 5. Operational model for the HDC-Manager Agent

The HDC-Manager constitutes an intelligent router that "brokers" discrete client-server interactions. However, tasks such as analyzing scientific data often require complex sequences of interactions to transport and manipulate data across distributed compute servers (e.g., signal processing and visualization tools). Defining such sequences across the relevant Agents can be difficult, particularly when those Agents participate in multiple task sequences. The Process Planner Manager supports a process-oriented distributed coordination model to solve this problem [9]. Interaction sequences are captured in process *scripts*. The Process-Planner coordinates the execution of such scripts by dispatching individual script elements as service requests to the HDC-Manager, acting as a driver to the HDC-Manager's router. Scripts promote maintainability and extensibility by centralizing the specification of sequences of distributed interactions. Moreover, the ability to combine Managers into hybrid control architectures illustrates the power of the Agent library's building block approach to designing distributed systems. A third Manager Agent, called a Server-Group, supports one-to-many client-server computing. This Agent decomposes complex requests into simpler services that independent server Agents can process concurrently. The Server-Group then collects and processes the results, assembling them into a single response. This interaction model is useful for groupware and decision support systems (e.g., collecting information dispersed across many databases and performing relational joins or analysis).

Symbiotics is working with NASA to refine and apply the NetWorks! Agent library to build several complex distributed systems. For example, NASA has developed various expert systems that "retrofit" the Launch Processing System for the Space Shuttle with automated FDIR capabilities. NetWorks! Gateway and HDC-Manager Agents are being utilized to integrate and coordinate these independent systems to work together cooperatively (e.g., by sharing data and diagnostic reasoning). A second NASA project is using Gateways to integrate tools for planning and scheduling ground operations for Shuttle missions. HDC-Manager and Process-Planner Agents are being used for script-based control of complex distributed decision support activities. Other contracts have investigated using library Agents: to design a framework for developing and managing software projects across heterogeneous computing platforms; to integrate tools for computer-aided design; and to develop group-based models for exploiting redundant applications to enhance reliability and corroborate problem solutions.

CONCLUSIONS

A central problem in developing distributed computing systems is enabling applications to interact across the network. Solutions to this problem are now being called "middleware," which encompasses all of the software between end-user applications and their networked host platforms. Other forms of middleware, such as network programming, RPCs, and conventional messaging tools, have one or more critical drawbacks, such as limited portability across heterogeneous platforms, intrusiveness, and limited reusability. NetWorks! middleware exploits object-oriented technologies for an advanced messaging solution that promotes modularity, maintainability, extensibility, and reusability. These characteristics are critically important to minimize development costs and to support lifecycle engineering of complex distributed systems.

NetWorks! provides non-intrusive peer-to-peer interaction models, which foster open, "plug and go" distributed computing architectures. Every application represents a potential server resource to other programs. The NetWorks! Agent library encourages an intuitive "building block" approach to developing distributed systems, in which complex distributed interactions are facilitated through service directories and simple process scripts. Complex distributed coordination functions are inherited from predefined Agent classes and customized to fit application requirements through high level APIs. NetWorks! Agents shield developers and end-users of distributed applications not only from network protocols, but also from data management functionality and complex control logic for handling messages. This technology thereby addresses the critical problems of accessing and coordinating data and computing resources in complex distributed systems.

ACKNOWLEDGMENTS

NetWorks! technologies have been developed with funding from the Small Business Innovative Research Program by the U.S. Army (contract DAAB10-87-C-0053) and NASA (contracts NAS10-11606, NAS10-11882, NAS5-31920, and NAS8-39343).

REFERENCES

- [1] G. Coulouris and J. Dollimore, *Distributed Systems: Concepts and Design*, Addison-Wesley, Reading, Massachusetts, 1988.
- [2] G. Andrews, "Paradigms for Process Interaction in Distributed Programs," *ACM Computing Surveys*, Vol. 21, No. 1, March 1991, pp. 49-90.
- [3] W. Stallings, P. Mockapetris, S. McLeod, and T. Michel, *Handbook of Computer-Communications Standards Vol. 3*, Macmillan, New York, 1988.
- [4] A. Birrell and B. Nelson, "Implementing Remote Procedure Calls," *ACM Transactions on Computer Systems*, Vol. 2, No. 1, February 1984, pp. 39-59.
- [5] D. Gifford and N. Glasser, "Remote Pipes and Procedures for Efficient Distributed Communication," *ACM Transactions on Computer Systems*, Vol. 6, No. 3, August 1988, pp. 258-283.
- [6] B. Meyer, *Object-Oriented Software Construction*, Prentice-Hall, New York, 1988.
- [7] R. Corbin, *The Art of Distributed Applications*, Springer-Verlag, New York, 1991.
- [8] R.M. Adler, "A Hierarchical Distributed Control Model for Coordinating Intelligent Systems," *Telematics and Informatics*, Vol. 8, No. 4, 1991, pp. 385-402.
- [9] R.M. Adler, "Coordinating Complex Problem-Solving Among Distributed Intelligent Agents," to appear in *Telematics and Informatics*, Vol. 9, No. 4, 1992.

516-60
150486

NO3-25577

THE DATABASE QUERY SUPPORT PROCESSOR (QSP)

Patrick K McCabe
Rome Laboratory
Rome, NY 13441

ABSTRACT

The number and diversity of databases available to users continues to increase dramatically. Currently, the trend is towards decentralized, client server architectures that (on the surface) are less expensive to acquire, operate and maintain than information architectures based on centralized, monolithic mainframes.

The database query support processor (QSP) effort evaluates the performance of a network level, heterogeneous database access capability. Air Force Material Command's Rome Laboratory has developed an approach, based on ANSI standard X3.138 - 1988, "The Information Resource Dictionary System (IRDS)" to seamless access to heterogeneous databases based on extensions to data dictionary technology.

To successfully query a decentralized information system users must know what data are available from which source, or have the knowledge and system privileges necessary to find out. Privacy and security considerations prohibit free and open access to every information system in every network. Even in completely open systems, time required to locate relevant data (in systems of any appreciable size) would be better spent analyzing the data, assuming the original question was not forgotten.

Extensions to data dictionary technology have the potential to more fully automate the search and retrieval for relevant data in a decentralized environment. Substantial amounts of time and money could be saved by not having to teach users what data resides in which systems and how to access each of those systems. Information describing data and how to get it could be removed from the application and placed in a dedicated repository where it belongs. The result: simplified applications that are less brittle and less expensive to build and maintain. Software technology providing the required functionality is off the shelf. The key difficulty is in defining the metadata required to support the process.

The database query support processor effort will provide quantitative data on the amount of effort required to implement an extended data dictionary at the network level, add new systems, adapt to changing user needs, and provide sound estimates on operations and maintenance costs and savings.

THE DATABASE QUERY SUPPORT PROCESSOR (QSP)

INTRODUCTION

The Database Query Support Processor (QSP) is the culmination of research and development that began with a particularly complex database conversion effort. In the early 1980's, Strategic Air Command (SAC) decided to migrate their entire intelligence support database to a completely different environment. Originally, SAC/IN was supported by a unique, home grown database management system developed specifically for SAC in the mid 1970's. In terms of maintainability this was intolerably expensive. To decrease maintenance costs, it was decided to migrate to a commercial product.

The database management system (DBMS) for the new system was the Cullinet DBMS. The Cullinet DBMS (called the Integrated Data Management System or IDMS) was considered by many to be the best DBMS at the time. IDMS was based on the network data model¹, which was consistent with SAC's existing data architecture.

Although the network data model was common to both databases, the hardware platforms and DBMS internals were completely different. The hardware platform in use was a Honeywell 6080; 4 CPU's, 1 MByte main memory (36 bit), and 3.8 GBytes (36 bit) disk storage. The target architecture was an IBM 3081; 4 CPU's, 32 MBytes (32 bit) main memory, 8.8 GBytes (32 bit) disk storage.

The conversion process was intensely manual. Software tools to assist this process were not available and had to be developed from scratch and on the fly. Change control procedures were lengthy and complicated. There were four distinct partitions constituting the development system at HQ SAC; one for development, one for integration, one for final testing, and a fourth for operational use. Physically moving the applications and data from one partition to the next was tedious. Many test errors were traced to missing pieces of software or incorrect versions of software modules being ported from one partition to the next.

Another requirement of the transition process was to provide simultaneous access to both systems. The sheer magnitude of the transition, with its inherently high technical risk, made a "knife switch" cutover approach an unacceptably high operational risk. The databases on both old and new systems had to be synchronized, and both systems required cognizance of what portions of the "operational configuration" were on which system. The existing user interface had to be maintained to the maximum extent possible. Users had to be insulated from the idiosyncrasies of each individual system².

NETWORK RESIDENT TRANSITION SUPPORT

The transition could have been orders of magnitude more difficult but for a unique element of SAC's architecture, the Micro-Programmable Controller (MPC). The MPC was an array of asynchronously operating microprocessors that shared a common backplane bus³. Developed originally to normalize the physical interfaces between quasi-intelligent workstations of various vendors and the Honeywell mainframe, the MPC evolved into a sophisticated distributed computing environment that was well ahead of its time.

Software was developed within the MPC to support simultaneous system access, minimizing changes to the user interface. Host resident software on the Honeywell system did not require modification and there was no need to develop throw away code on the IBM system. Software implemented on the MPC was essentially an

1 Most aspects of data models are extremely well covered in [MAR77]. Cullinet was absorbed by Computer Associates in the late eighties.

2 Additional information on the transition effort is provided in [RAD85].

3 The MPC predated general acceptance of local area networks. It still provides some network services, but has mostly been supplanted by a local area network. The local area network consists of clusters of IEEE 802.3 LANS connected by an FDDI backbone. Additional details pertaining to the MPC may be found in [RAD86].

extension of the network support functions already provided. Unfortunately, this software was essentially thrown away since it would have no purpose once the transition phase was complete.

The difficulties encountered during the transition effort made it clear that automated tools were required for future database transitions. It was also clear that simultaneous access to multiple databases would be a required capability for future systems. The network itself was the logical provider of these capabilities. What exactly these services should be and how the network should provide them was the primary question. Some sort of dictionary/directory would be required that provided database access support services, but what was required beyond that wasn't clear.

DESIGN CONCEPTS FOR DATABASE UTILITIES

As a result, a study effort entitled "Design Concepts for Database Utilities" was initiated to better define the characteristics of network level database access utilities. An architecture for an "Integrated Data Network (IDN)" was developed⁴. The architecture consisted of a three level hierarchy of six types of processors, four of which were specific to the IDN (see figure 1.)

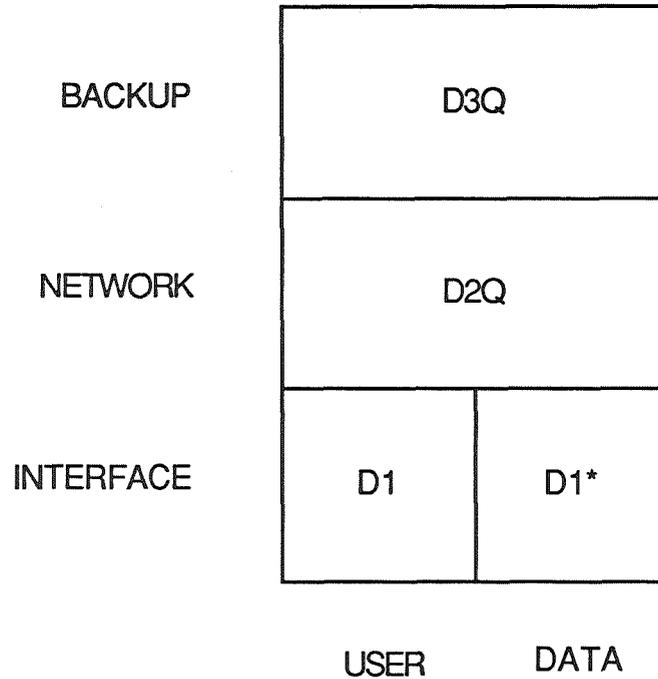


Figure 1
Hierarchy of Processors

The user node corresponds to the processor at which the application or user requesting data resides. Data nodes are the physical repositories of the requested data. User nodes and data nodes are considered outside the scope of the IDN.

At the interface level of the IDN architecture are D1 nodes and D1* nodes. D1 type nodes interface user nodes to the network, accept queries, perform first order validation of the queries, and assemble query responses. D1* type nodes interface data nodes to the network, receive subqueries directed to specific data nodes, accept responses from the data nodes, and compose aggregate responses for transmission to D1 nodes.

At the network level of the IDN architecture are the D2Q nodes. The D2Q nodes complete query validation, dispatch subqueries, and control query execution. These nodes are core to the IDN architectural concept, providing the actual dictionary, directory, and query support services required.

The D3Q node is at the backup level and serves to provide backup facilities for all other types of node, except the user node. Additionally, contents of data nodes can be replicated on D3Q nodes. Replicating data (in the long haul network environment) can improve performance by balancing communication load and supporting fault

⁴ See [RAD86.1] for more details. It is also important to realize that the context of this effort was a wide area (if not global) information network. Performance and fault tolerance were critical design considerations.

tolerant operations. Data node failures won't halt query activity. The resultant network architecture is depicted in figure 2, below⁵.

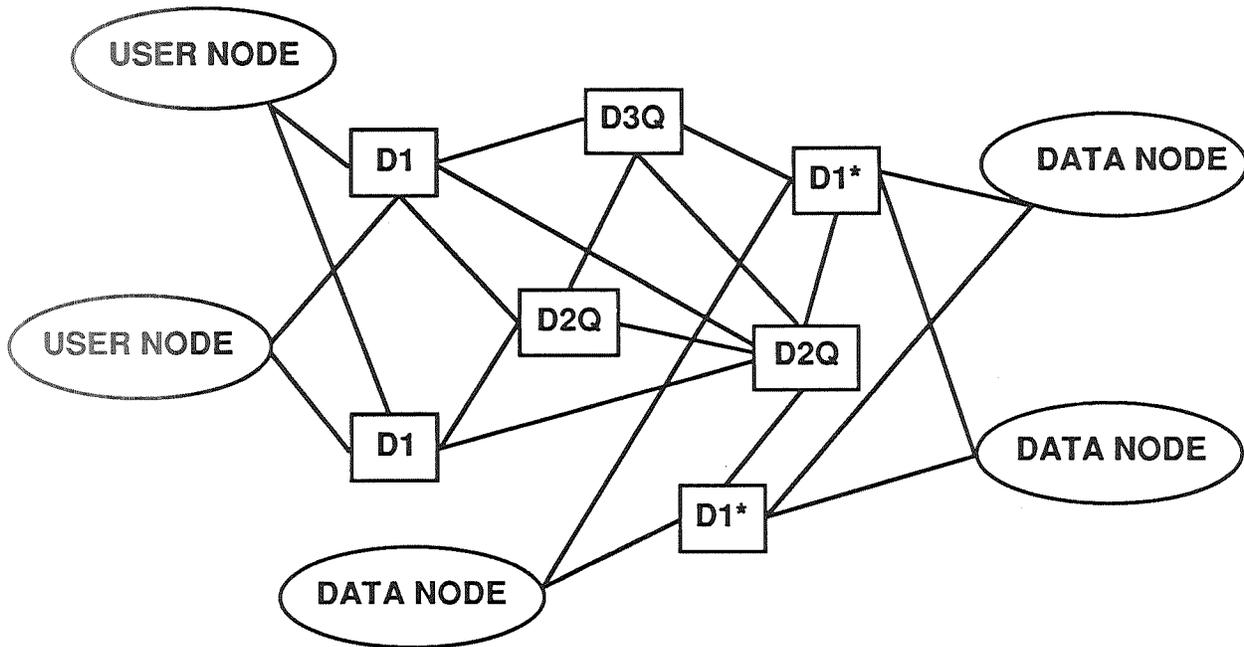


Figure 2
IDN Architecture

⁵ The architecture was developed deliberately to maximize functional redundancy. The figure illustrates this concept by showing multiple paths between each node. At least one of these redundant paths connects to a shadow node, a node capable of acting as a hot backup for a similar node.

DATABASE QUERY SUPPORT PROCESSOR

During the effort it was realized that the same technology applied to local area networks as well. Implementation details would differ due to differing bandwidth, topological, and fault recovery characteristics of wide area networks versus local area networks. Within the local area network environment, the functionality of the D1 node would be absorbed by the user's workstation, the functionality of the D1* node would be absorbed by the data node, and the D3Q node would constitute the QSP. Since the D3Q provides all the functionality of the D2Q, with the addition of replicated data from selected data nodes, the D2Q can be eliminated as a separate device (see figure 3, Notional QSP Architecture).

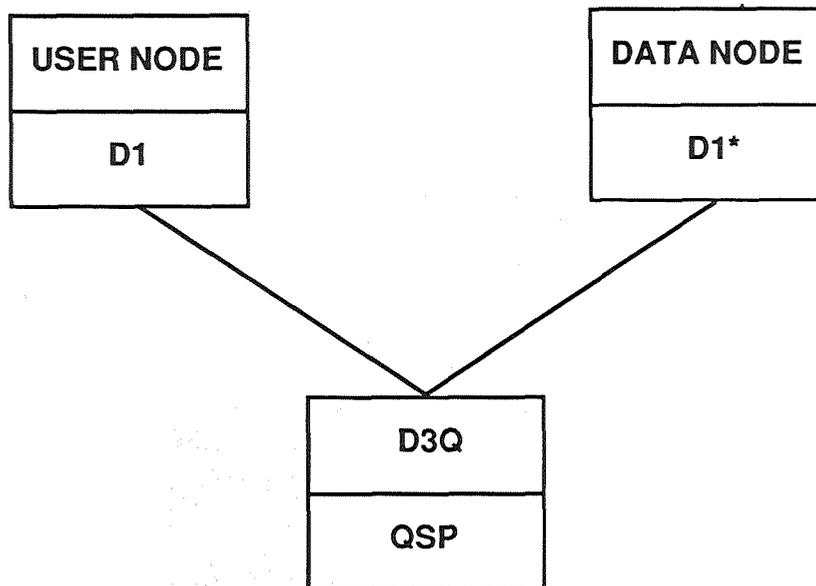


Figure 3
Notional QSP Architecture

For the proposed solution to be effective, it had to have the characteristics of an active, in-line data dictionary at the network level⁶. This meant that all activity against the databases in the network, including application development, database modification and maintenance, and routine database access had to utilize services provided by the utility. The methodology for operation of the IDN and subsequently the QSP, was based on the emerging Information Resource Dictionary System (IRDS) standard⁷. In other words, the functionality of the D2Q or D3Q nodes discussed above, was based on the IRDS standard.

THE INFORMATION RESOURCE DICTIONARY STANDARD

The motivation for the development of the IRDS standard, ANSI X3.138 - 1988, was the proliferation of redundant and inconsistent data. The data dictionary system was seen as a key tool for the effective management of information resources and reduction of inconsistent, redundant data. A number of incompatible, stand alone data

⁶ Detailed discussion of the philosophy behind data dictionaries and their characteristics is provided in [ROS81].

⁷ There were two efforts initiated about the same time to develop standards in this area. The American National Standards Committee for Information Systems (X3) began work on a standard for an "Information Resource Dictionary System." The National Institute of Standards and Technology (NIST, formerly the National Bureau of Standards) effort focused on the development of a Federal Information Processing Standard for Data Dictionary Systems. Both groups had identical goals and similar approaches [QED85]. Both efforts were merged in 1983 and the result was the IRDS [ANS88].

dictionary systems were on the market, and each database management system had closed, internal implementations of data dictionaries (if they had any). It was perceived as necessary to develop a standard for data dictionary software⁸.

The IRDS standard describes a four level information architecture, level 2 and 3 of which constitute Federal Information Processing Standard (FIPS) 156 (see figure 4, IRDS Architecture). Each level describes and controls the lower level. The first level, Information Resources, is the data in your database. The standard does not apply at this level, although it must accommodate it. The second level, the Information Resource Dictionary (IRD), is the data in the data dictionary, which describe the data in the database. One likely extension of the IRDS approach is to extend the control function from level 2 to level 1. As you might expect, the data dictionary is itself a database that consists of data elements and relationships. Definitions of the data elements and relationships that constitute the data dictionary must be managed. The third level of the IRDS Standard, the Information Resource Dictionary Schema, consists of the definitions of the data elements and relationships contained by the data dictionary. The fourth layer is called the Information Resource Dictionary Schema Description, and consists of data that describes the IRD Schema (level 3).

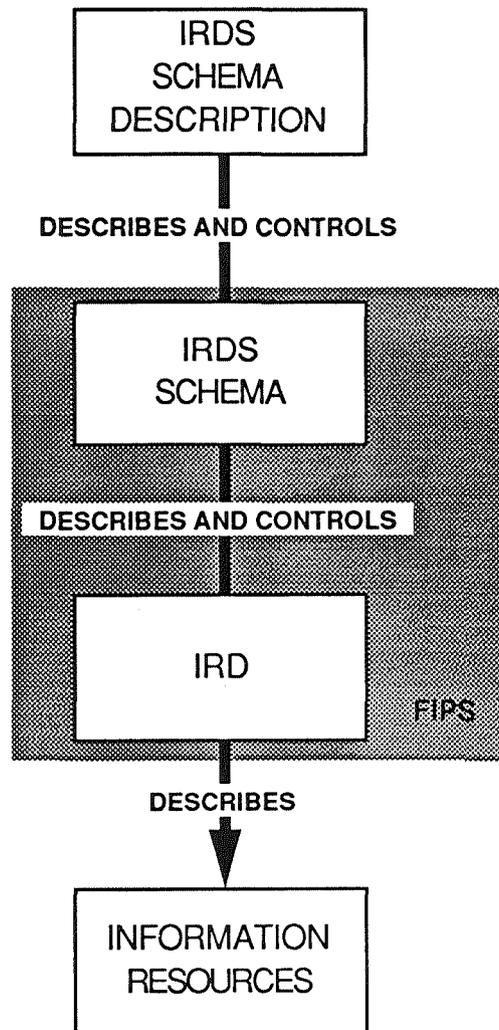


Figure 4
IRDS Architecture

⁸ See [QED85]. The standards committee took the approach that the standard should specify the characteristics of an interface to a data dictionary and the functionality that a data dictionary should provide. They wisely avoided the mistake of trying to dictate how to implement the dictionary itself.

Key to the concept of levels of description is the corollary that the higher the level, the simpler the model required to describe it. What is left is a mechanism that anyone can use to retrieve data relevant to a specific query. Services provided at each level take care of details such as how to determine what data is available, how to locate it, how to request it, how to navigate the database to get it, and how to put it together into a usable product.

The results of the Design Concepts for Database Utilities work were used by the performing contractor to develop a commercial product in this area. They were successful in obtaining SBIR phase I and phase II funding, and did build a prototype⁹. Rome Laboratory became aware at this time that several organizations were working on similar capabilities.

By 1990 it was apparent that the technology required to support network level database support utilities was mature. The last set of questions requiring answers prior to operational implementation of the technology pertained to performance and policy. More specifically, how much overhead would be introduced into operational systems to achieve what degree of benefit (in terms of flexibility, operations and maintenance savings, etc.). Additionally, simultaneous access to multiple databases adds a new dimension to security policies and procedures, which must be fully understood before implementation.

QSP STATUS

In 1991, the Database Query Support Processor (QSP) effort was initiated to answer these questions. The effort presupposes the availability of network level database support systems with the following capabilities;

- a. To retrieve data from multiple databases irregardless of data location, database architecture, or database navigation constraints.
- b. To support the definition, modification, administration, and maintenance of:
 - (1) A network level schema describing the totality of information available from all databases in the network.
 - (2) Network level subschemas, which are logical subsets of the network level schema and assigned to specific classes of operational users.
- c. Provide tools to assist database administrators in defining specific database views for inclusion in the network level schema.
- d. Manage and control the definitions of, inter-relationships among, and definitions of inter-relationships of: data elements, data structures, applications, products, user descriptions, and information requirements.

During 1992 and 1993, the QSP effort will focus on collecting quantitative data such as:

- a. Volume, patterns, and types of network traffic generated by the QSP.
- b. Volume, patterns, and types of network accesses to the QSP.
- c. Elapsed time from issuance of a query at a workstation to its receipt by the QSP.
- d. Elapsed time from receipt of a query by the QSP to generation of all subqueries.
- e. Elapsed time from subquery generation to subquery issuance by the QSP.
- f. Elapsed time from issuance of a subquery to receipt by host resident QSP interface software.
- g. Elapsed time from issuance of data request by the host resident QSP interface software to that software's receipt of the host's response.
- h. Elapsed time from issuance of subquery response by the host resident QSP interface software to receipt of the response by the QSP.
- i. Elapsed time from receipt of all subquery responses to the issuance of a query response by the QSP.
- j. Elapsed time from issuance of query response by the QSP to receipt of the response by the workstation.

⁹ The Small Business Innovative Research (SBIR) program provides up to \$50,000 for phase I efforts and results in a specification for phase II implementation. Phase II provides up to \$500,000 for implementation of the idea. Phase III is usually contractor funded and results in a commercial product (with some limited Government rights). See [RAD90] and [RAD90.1] for more information on the SBIR efforts.

The effort will wrap up in 1993 with a comprehensive analysis of collected data in the context of an operational environment. Implications to security policy and accreditation, hardware and software short comings, and operations and maintenance costs will be assessed. Flexibility of the QSP approach will be assessed with respect to the amount of work required to accommodate new databases, changes to old databases, and to initially implement the QSP in an operational network. This data will be used to build a specification for a production version of the QSP.

CONCLUSION

The benefit of the QSP is in the network level support services made possible by the active, in-line repository at the heart of the device. Knowing the relationships among data elements and applications across system boundaries allows better control over change. The ripple effect induced by modifying data elements or applications can be identified in advance and more effectively priced. Additionally, data elements may already exist somewhere in the network that meet the needs of a proposed development, minimizing new development.

Additional benefits could result from adding system documentation to the information available in the network. From the QSP's perspective, documentation can be treated as just another database. Network level information pertaining to relationships among documentation, data elements, applications and other elements of the information environment could be maintained in the QSP. This capability makes update of relevant system documentation an integral part of application or database development, rather than an afterthought.

Data element and application standardization are also supported by the information contained in the QSP repository. The information necessary is already available, all that would remain is to define the rules. Triggers or other mechanisms provide the vehicles for implementation.

The QSP effort will provide hard data on which to base future implementation decisions. Specifically, which services to implement and to what extent to implement those services in operational IDHS systems. Start up costs and the operations and maintenance tail required will also be determined. In the long run, the QSP should provide real benefits in terms of more flexible and robust information systems, with lower operations and maintenance costs.

BIBLIOGRAPHY

- [ANS88] American National Standard Information Resource Dictionary System, ANSI X3.138-1988 (Federal Information Processing Standard 156), ANSI, New York, 1988.
- [MAR77] Martin, James, Principles of Data-Base Organization, 2nd ed.;Prentice Hall, 1977.
- [QED85] AOG Systems Corporation, The Draft Proposed American National Standard Information Resource Dictionary System, QED Information Sciences Inc, Wellesley, MA, 1985.
- [RAD85] Planning Research Corporation, Anderson, Richard D and Gerald H. Paes, SAC IDHS Improvement Project, RADC-TR-85-187, October 1985, Secret.
- [RAD86] McCabe, Patrick K., Michael J. Wessing, and Ernst K. Walge, SACINTNET Future Study, RADC-TR-86-144, August 1986.
- [RAD86.1] AOG Systems Corporation, Design Concepts for Database Utilities, RADC-TR-86-48, April 1986.
- [RAD90] AOG Systems Corporation, Local Area Network Schema Server, RADC-TR-90-376, December 1990.
- [RAD90.1] AOG Systems Corporation, Local Area Network Schema Server, RADC-TR-90-375, December 1990.
- [ROS81] Ross, R.G., Data Dictionaries and Data Administration, New York, AMACOM, 1981.

ORIGINAL PAGE IS
OF POOR QUALITY



omit

**MANUFACTURING TECHNOLOGY
PART 1**

PRECEDING PAGE BLANK NOT FILMED

167.

166
~~INTENTIONALLY BLANK~~

THE UNIVERSITY OF CHICAGO

517-35
150487
p-8

APPLICATION OF AN ON-MACHINE GAGE FOR DIAMETER MEASUREMENTS

Kevin G. Harding
Industrial Technology Institute
P.O. Box 1485
Ann Arbor, MI 48106

ABSTRACT

This paper describes the design analysis and application of a laser based gage made specifically for measuring parts on the machine tool to a high accuracy. The tri-beam gage uses three beams of light to measure the local curvature of the part in a manner similar to a V-block gage. The properties of this design include: calibration that is independent of the machine tool scales, non-contact damage free operation, low cost of the gage, and the ability to measure parts in motion.

INTRODUCTION

Increasing tolerances in machining have driven the need for closer and closer control of the process. To gain the greatest performance out of a machine tool, it is necessary to have the dimensional data available for minute adjustments to be made to machine offsets. In the past, many of these dimensional checks have been performed off-line, often some time after a part has been made. The time lag from part manufacture to measurement has meant that a drift or needed adjustment would not be made until many parts had already been manufactured at the previous settings. The most desirable place to make a measurement is right on the machine tool. Any measurement on the machine tool must be made with minimal affect on the cycle time of the machining least a cost be added on to the manufacturing cycle.

On-Machine Touch Probes

To facilitate on-machine gaging, many machine tool builders now supply a touch probe option. The touch probe is used either as a tool (put in place of the tool when measurements are made), or is otherwise attached to the machine. A measurement is typically made by first touching the touch-trigger probe to a reference datum on the machine tool. The probe is then moved by the machine mechanism to the points on the part to be measured.

The difficulties associated with making on-machine measurements with a touch probe are well known. In turning operations, the touch probe must measure the part along the radius of the part. To the extent that the measurement axis defined by the probe deviates from a line along the part radius, there will be an error in the measurement. That is, if the measurement line of the probe is off-set from the radius of the part, the probe will actually measure along a cord of the part, rather than the full radius. Once the part is sensed, the size and extend of the probe tip must be accounted for as an offset in the measurement. If the normal to the part is not know, this offset correction can be difficult.

Making a probe measurement with the machine tool system also has the problem that the part dimension is checked with the same mechanism which makes the part. If the scales on the machine tool are off, then the measurement may be wrong, to the same degree, for both making and measuring the part. Therefore,

using the machine tool to make the measurement of the part does not provide an independent check of the machine accuracy.

A final difficulty with making measurements with an on-machine touch probe is the time involved. As the machine itself is used in the measurement, the time spent making the measurement adds to the effective machining time. Typically, the part must be stopped before making measurements. If multiple measures around the part is desired (as is sometimes the case), the part must be rotated to the those positions to be measured and held for the time it takes to approach the part with the touch probe to make the measurement. The touch probe can not be moved into the part quickly, as it can easily be damaged in collision with the part. The measurement speed of a touch probe type system is therefore typically limited to 1 to 2 points per second at best.

The problems described do not prevent the measurements from being made, but they do have three primary effects:

the cycle time for the machine tool utilization is increased by the measurement time, which also limits the measurements made,

the measurement is checked only against the machine tool itself, thereby not necessarily finding any error in the machine tool positioning that may have machined a part wrong,

down time for the gage can be a problem because of damage of the touch-trigger probe.

These effects have been a limiting factor in supplying reliable, on machine gaging for offset corrections.

Off-Machine Gage Options

The alternative to on-machine gaging is to perform near machine gaging using such tools as machine vision, mechanical calipers, or laser micrometers.¹⁻⁷ Machine vision has been used to the accuracies desired (about 2.5 micron) by means of optically based versions of coordinate measurement machines. In these systems, a camera with a small field of view is moved across the field by a precision encoded stage system, which may be a gantry very similar to a traditional coordinate measurement machine (CMM). An off line measure, by even the faster optical CMM systems is still removed from the machine, and hence there is a temporal lag between the part completion and the availability of the measurement data, as described before.

The laser micrometer field is one which has become well established for near machine gaging. Laser micrometers offer the advantage of being a noncontact method, not prone to physical damage due to contact with heavy duty machines as many mechanical caliper systems suffer. A laser micrometer obtains it's measure by scanning or just shadowing a beam of light around the part, to create a silhouette of the part diameter. In the scanner based laser micrometers, the diameter is determined by the time it takes the beam to pass the object (the time during the scan for which the laser beam is shadowed by the part).

Laser micrometer systems that use a simple shadow use a static sheet of collimated laser light which is shadowed onto a linear detector array. The edge of the shadow of the part has a distinct shadow shape to it cause by the diffraction of the laser light passing the part. The linear array can provide a sampling across the shadow edge diffraction pattern which allows the edge location to be determined to much less than a pixel. Typical accuracies for either type of laser micrometer are a part in 10,000 to a part in 20,000. Particularly with the laser scanning systems, diameter measurement resolutions of better that 0.25 microns (10 millionths of an

inch) can be made though the use of a large number of measurement samples (one hundred measurements can be made in a second or less with these systems). A large degree of environmental stability is needed for measurements of less than a micron to be meaningful.

The laser mike has the advantage over a single touch probe in that it necessarily measures the greatest diameter of the part (as it is looking along parallel tangents on opposite sides of the part), and is noncontact, thus avoiding damage. As a noncontact probe, the laser micrometer can make measurements very quickly (a few hundred per second) as there is not danger of collision of the light beam with the part under test. In addition, the laser micrometer does not rely on the machine tool accuracy, but rather makes an independent measure of part diameters. Laser micrometers are currently established in industry for such applications as measuring wire and extruded tubing in-process. The measurement for wire production applications is made with the measured material in motion. Any multiple sampling of a moving product, of course, will produce an average measurement along some length of the product, but that is quite acceptable in this type of application.

For the application of the laser micrometer to machine tool operation, there are some drawbacks. Because the laser micrometer must surround the part (with instrumentation on both sides of the part), for an accurate measure, it is typically not practical to consider putting such a device in the machine tool itself. The long, open light path of the laser micrometer can also be a problem in dirty environments as air turbulence, heat, or other airborne material can deviate the light beam (as a function of how far the light must travel) and produce bad data. Maintaining the stability of the transmitter on one side to the receiver on the other side is an important structural and environmental concern which must be dealt with in using a laser micrometer. In typical applications where the laser micrometer can be rigidly mounted, the mechanical stability problems have been well addressed by the commercial vendors of this equipment (air environment has remained a problem). However, if the gage is to be moved, such as to measure a turned part in the chuck, these stability problems could be very limiting.

Therefore, there remains a need for a more effective gage for on-machine gaging of parts. The desirable features for such a gage for application of outer diameter turned part gaging would include:

- high speed measurement capability,
- accuracies of 2 to 5 microns,
- only single side access required,
- accuracy independent of machine scales,
- high damage resistance for in-machine use,
- limited sensitivity to heat and air contaminants of machine tool environment,
- ability to measure moving parts (such as during spin down) to minimize time cost of the measurement.

The application we will be addressing is limited to outer diameter (OD) measurements.

TECHNICAL APPROACH

The tri-beam gage works on the same principle as the standard v-block gage. A v-block gage consists of a physical v-block, in which the cylinder to measure is placed. The line contacts of the v-block with the cylinder establishes two tangent points on the cylinder. A micrometer is mounted in the apex of the v. The micrometer is advanced till it just contacts the cylinder (this can be difficult to insure). The reading from the micrometer then permits a calculation of the cylinder curvature in that region, and hence a measure of the cylinder diameter.

Mechanical v-block gages are as well established as the caliper. V-block gages have the desired property that the measure is not dependent on any outside positioning of the gage. Mechanical v-block gages have typically be used on larger cylindrical objects where surrounding the part with a caliper or micrometer would require and impractiably large gage that may be difficult to manage. Therefore, the v-block gage only requires a small part of the surface of the cylinder to make the measure, and does not require surrounding the part (it need access only a single side). Examples of where this type v-block gage has been used includes large gun cylinders, trees, and optical components (where, of course, on a lens only part of the curve may exist).

The tri-beam gage is simply the optical equivalent of the v-block gage. As shown in Figure 1, in place of the V-structure we use two beams of light, detected by linear array detectors. In place of the vertex mounted micrometer we use a sheet of light and a detector array, similar to that used by commercial laser micrometers.

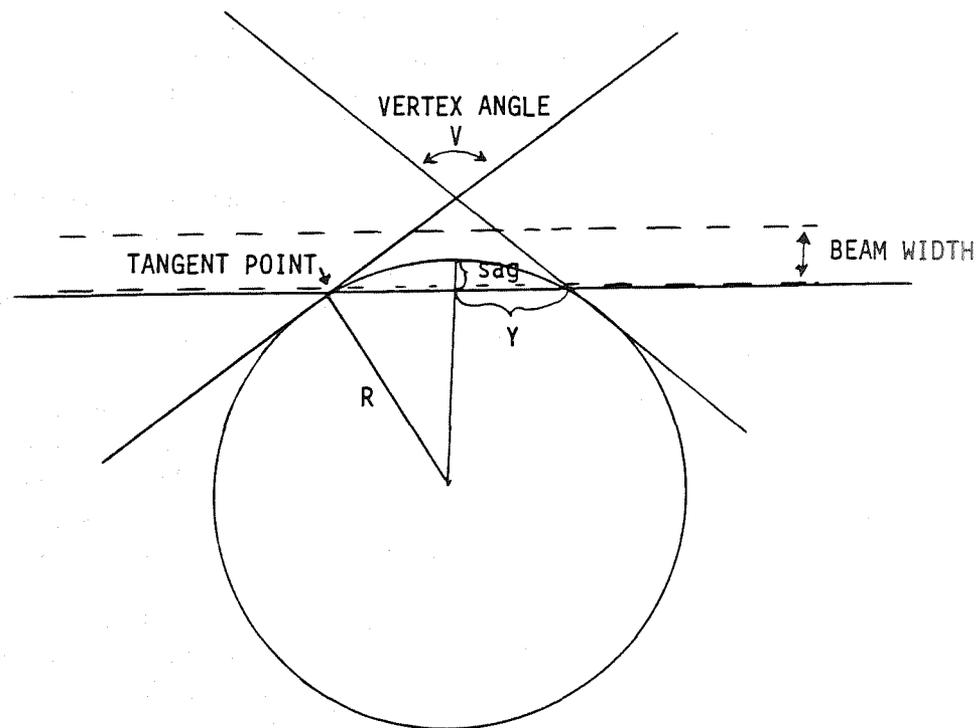


Figure 1. Diagram of the tri-beam gage concept.

For the purpose of establishing the two tangent points, the tri-beam gage simply shadows the light onto the two side detectors. The measurement range of motion of the position of the cylinder edge need not be very large, as the purpose of the two side detectors is simply to establish the two tangent points at the time we make the center measurement which will be equivalent to the vertex micrometer. The third beam used to replace the vertex mounted micrometer must have the accuracy, and range of geometric measurement to accommodate the various diameters of interest for the particular gage.

There are a number of parameters which affect the performance of the tri-beam gage. The primary parameters include gage geometry, and gage stability in use. The gage mechanical stability is driven by the mounting of the three detectors and light sources with respect to each other. By mounting these components rigidly in close proximity to each other, the effects of vibrations or thermal expansion of the unit can be minimized. For example, mounting the components with a linear offset from each other of 5 centimeters on steel

would produce a drift of 0.5 microns (20 microinches) per degree C. However, mounting these components in a symmetric configuration in a 30 degree cross configuration would reduce this drift to 0.2 microns for a temperature differential across the 5 centimeters of a degree centigrade. Therefore, it becomes only necessary to compensate for uneven heating by means of controlled radiative heating to maintain negligible mechanical drift with temperature.

The other consideration relating to gage geometry is the actual angles between the beams of light used to form the optical V-block. As the V configuration becomes shallower, a greater accuracy is needed in the measurement at the vertex point to maintain the same accuracy on the diameter measurement, and the greater the significance of any errors in the optical V measurements. For example, a standard array sensor could be used to measure a radius range of 50 millimeters (diameter range of 100 millimeters or 4 inches) with an angle of the V of 60 degrees. The detector in this case would need to measure to an accuracy of twice as good as the desired radius measure (4 time the diameter). The same sensor could measure a radius change of 100 millimeters (diameter range of 200 millimeters or 8 inches) with a V angle of 120 degrees, but would need to be 4 times more accurate than the desired radius measurement accuracy (8 times better than the diameter). Therefore, to obtain 2.5 micron accuracy, the sensor would need to be accurate to about 0.3 microns. This accuracy is within the limits of state-of-the-art detector systems.⁸⁻¹⁰

As with any device of this type, the part needs to be clean to make a reliable measurement. If any chips or other debris are present, they may be measured along with the part to produce an erroneous measurement. This need for some degree of cleanliness is actually true for any diameter measurement (even a chip in a mechanical V-block will lead to a wrong measurement). Cleaning of the part can easily be accomplished by blowing off the surface. Such a provision can be built right into the gage head itself.

Coolants and other debris in the air is not of major concern with the tri-beam gage concept. As the gage detects edges, rather than a light level to determine diameter, as long as there is some light getting through, the measurement is possible to make. If there is noticeable debris present on the gage, this may contribute to the measurement noise. A standard way of dealing with contaminants with this type of gage is to use a regular air flow or "air curtain" to keep any contaminants moving, and therefore preventing any buildup.

The performance objective for this gage were based upon input from machine tool builders and users. We found that large numbers of data point are not needed because the controllers can not use the extra data anyway. Most on-machine measurements are made to provide just simple offsets from some known dimension. The speed fo the tri-beam gage may provide additional data relating to ovality, and permit multiple sampling locations in a short period of time. The performance objective derived are summarized as follows:

Table 1. PERFORMANCE OBJECTIVES

- Accuracy: 1-2 micron (0.00004 inches)
- Speed: 10 measures per second or better
- Environment: \pm 30 degrees C
- Part Sizes: 5 cm to 15 cm (2 to 6 inches)
- Calibration: internal
- Environmental: chips, coolant (between cuts)
- Price: under \$10,000

The system we built uses the simplest configuration of optics and mechanics. Based upon a 60 degree vertex angle, we were able to design the system to cover a range of diameters of 2 to 6 inches (5 to 15 centimeters). A diagram of this system is shown in Figure 2. The laser light is collimated and is shadowed

directly onto the detectors. The interference pattern produced is shown in Figure 3. By matching to the well defined pattern, the position of the edge can be found very precisely. This allows for a measurement resolution of a few microns over a range of 4 inches of diameter. A picture of the sensor is shown in Figure 4. The structure was made very solid to minimize the effects of vibrations, heating, and minimize the chance of damage when operating on the tool carousel.

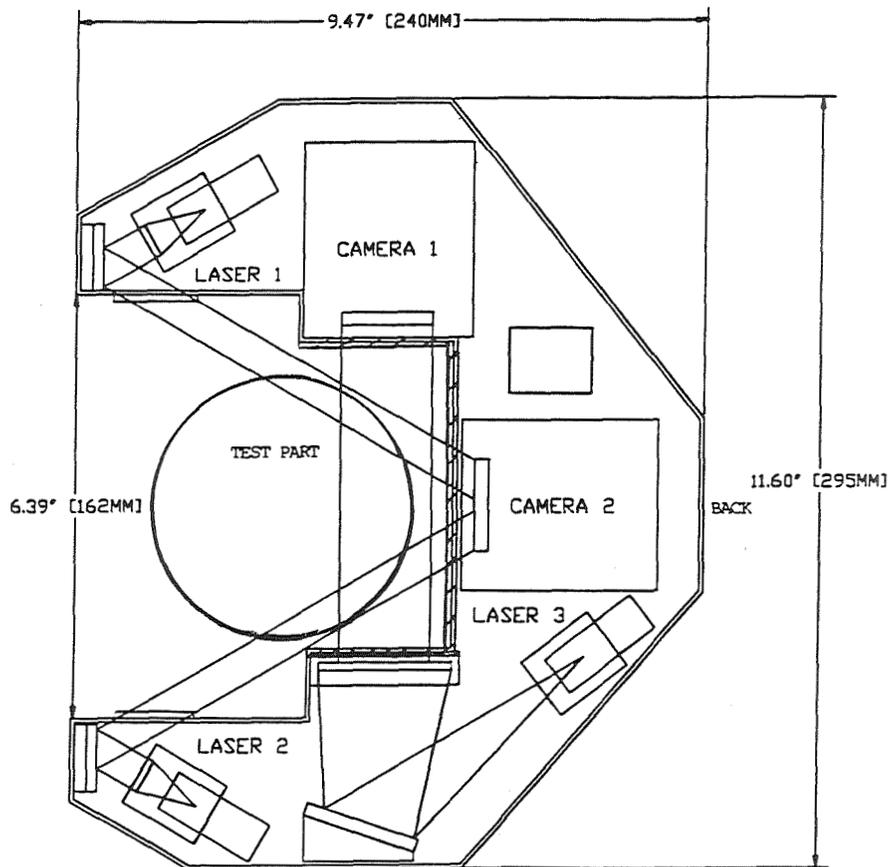


Figure 2. Diagram of 4 inch range tri-beam gage made by ITI.

APPLICATIONS AND CONCLUSIONS

The primary application for this gage is as a means to obtain immediate feedback to the machining process. When measurements are taken after the fact, compensations may be made which over correct the offset of the tool, leading to a continued swing plus and minus of the ideal. With information available while the part is on the machine, potentially subsequent cuts can be adjusted to correct errors on that particular part.

In general, the primary application motives for this gage are as summarized below:

- Provide On-Machine Gaging for Immediate Feedback
- Compensate Subsequent Cuts
- Monitor Machine/Tool Condition

- Obtain a Measure Independent of Machine Scales
 - Find Centering Errors
 - Monitor Machine Drifts

- Minimize Machine Time Interruption of Measurements
 - Measure During Spin-down
 - Provide High Measurement Speed

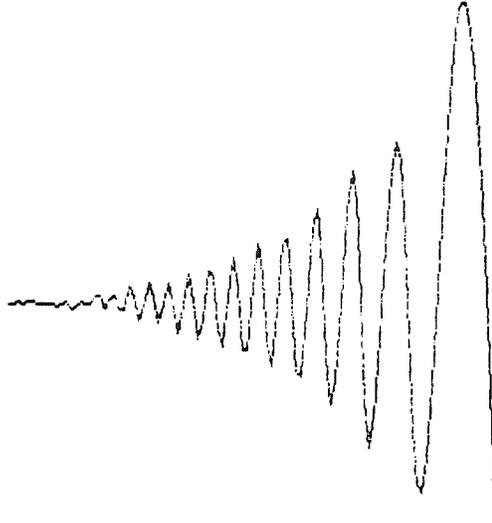


Figure 3. Interference pattern from the part edge used to make the measurements.

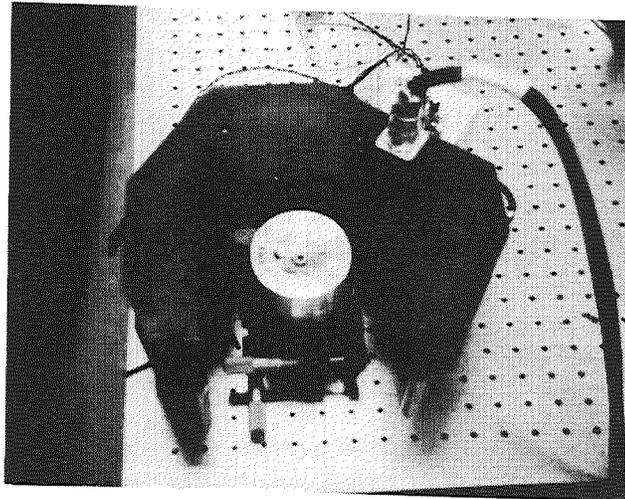


Figure 4. Picture of the prototype tri-beam gage made for on-machine tests.

As with any new gage, the ultimate applications will depend on the performance and final operational parameters of the gage system. With increasing abilities for machine controllers to use real-time data, on machine gaging such as this is expected to become an important tool for producing quality parts in the future.

ACKNOWLEDGEMENT

This work was funded by the US Air Force Manufacturing Technology Directorate under contract number F33615-91-C-5704, for whose support we are grateful.

REFERENCES

1. Tomohisa Mikami, Fumitaka Abe, Tohru Satoh and Tadashi Matsuda, "Ultra-High Precision UV Laser Scanning System," ICALEO Proceedings: Inspection, Measurement, and Control, Vol. 39, 1983.
2. J.C. Erdmann, R.I. Gellert, and R.L. Skaugset, "Laser Gage for Measuring Changes In The Surface of a Moving Part," Applied Optics, Vol.19 No. 18, 1980. patent # 4,148,587
3. Donald W. Whitney, "Optical Gaging-Crankshaft," Vision 85 Proceedings: Process Control and Gauging, pp. 7-1 to 7-11, 1985.
4. James A. Soobitsky, "Scanning Laser Diameter Gages for Industrial Use," SPIE Proceedings: Opto-mechanical and Electro-Optical Design of Industrial Systems, Vol. 959, pp. 193-209, 1988.
5. Henry C. Turko, "Industrial Gaging Technique With Emphasis On Scanning Laser Beams," ICALEO Proceedings: Inspection, Measurement, and Control, Vol. 33 pp. 1-10, 1982.
6. Larry P. Norman, "Optical Dimensional Gaging And It'S Integration To The Manufacturing," Vision 89 Proceedings: Vision for Manufacturing Cells, pp. 11-31 to 11-38, 1989.
7. M. Nohda, "Method Of And Apparatus For Measuring Radius," Applied Optics, Vol. 25 No. 15, 1986. patent # 4,572,628.
8. Diana Nyssonen, "Practical Method For Edge Detection And Focusing For Linewidth Measurement on Wafers," Optical Engineering, Vol. 26 No. 1, pp. 81-85, 1987.
9. Oleh Tretiak and Guo Yao Yu, "Curve-Fitting Method For The Measurement Of The Resolution Of digital Image Devices," Optical Engineering, Vol. 25 No. 12, pp. 1312-1315, 1986.
10. R.C. Anderson and C.S. Anderson, "Signal Processing Using Only Fourier Phase," Optical Engineering, Vol. 25 No. 12, pp. 1316-1319, 1986.

518-31

150488

P-4

N93-20579

ON MACHINE CAPACITANCE DIMENSIONAL AND SURFACE PROFILE MEASUREMENT SYSTEM

Ralph Resnick
Extrude Hone Corporation
8075 Pennsylvania Avenue
Irwin, PA 15642

ABSTRACT

A program was awarded under the Air Force Machine Tool Sensor Improvements Program Research and Development Announcement to develop and demonstrate the use of a Capacitance Sensor System including Capacitive Non-Contact Analog Probe and a Capacitive Array Dimensional Measurement System to check the dimensions of complex shapes and contours on a machine tool or in an automated inspection cell. The manufacturing of complex shapes and contours and the subsequent verification of those manufactured shapes is fundamental and widespread throughout industry. The critical profile of a gear tooth; the overall shape of a graphite EDM electrode; the contour of a turbine blade in a jet engine; and countless other components in varied applications possess complex shapes that require detailed and complex inspection procedures. Current inspection methods for complex shapes and contours are expensive, time-consuming, and labor intensive.

INTRODUCTION

An effort entitled "On-Machine Capacitance Dimensional and Surface Profile Measurement System" funded by the U.S. Air Force Wright Patterson Laboratories PRDA Program and performed by the MetreX Division of Extrude Hone Corporation seeks to address the needs of complex shape measurement and improve upon the use of the present technologies for performing those measurements. Recent advances in capacitance dimensional and surface roughness measurement provide a potential methodology for the verification of complex shapes. By building a capacitance probe or array and then scanning the workpiece, valuable information can be obtained about the surface under the sensor(s). This data can then be assimilated, translated, and transferred to the machining process controller for setup or corrective action. Since the cost per capacitance sensor is low and processing time is short, full form verification could be done quickly and inexpensively. Careful design and integration would permit the new measurement system to be installed in or near process.

The objective of this program is to develop and demonstrate a system to check and quantify dimensional information about complex workpiece shapes in a process environment on an automated machine tool. The proposer has established a partnership with Carrier Corporation of United Technologies on this program to establish key requirements, to review approaches and programs, to provide components, and to test results in a present process. Researchers at the University of Washington are working on design and analyzing mathematical models of the system as well as performing the environmental testing. The National Institute of Standards and Technology is providing the temperature, stability and dynamic testing and defined interface standards based on recently implemented Dimensional Measurement Inspection Specifications (DMIS).

PURPOSE AND OBJECTIVE

Identification of the Problem

The verification and inspection of complex shapes is far-reaching and significant. An improper profile on a gear can cause misalignment increasing gear noise and decreasing gear life. A cam with the wrong contour can cause erratic engine performance. The aircraft engine industry has many important complex shapes on important components and is highly representative of the requirements of industry in general. Aircraft turbine engine components require that their geometric form be accurate and consistent.

More efficient and predictable engine performance can be established if the engine components are maintained within the design criteria. To guarantee the components will meet performance specifications requires that the specific engine parts be gauged to the design criteria. For some simple components this gauging task is straightforward and easily performed; however, on complex components with intricate shapes and contours, specifically turbine blades and disks, Integrally Bladed Rotors (IBR's), and impellers, the gauging and design verification is detailed and complicated. Since the production of these components is performed on complicated machines with complex parts programs or intricate processes, the feedback of this test information is vital and timely for process and program control.

Present techniques for form measurement of complex aircraft engine components include Coordinate Measurement Machines (CMM's), Light Sectioning, SigmaFlex gauges, and even templates and feeler gages. All of these techniques are complex and/or labor intensive. The procedures are performed off-machine, at times even off-site, complicating the ability to transfer the valuable correction data to the manufacturing cell. Fixturing, programming and setup requirements are extremely time consuming and expensive. It has been estimated that the post-process dimensioning of a workpiece roughly accounts for between 20 to 40 percent of the total time to complete a machining operation, depending upon the complexity of the cut and the number of tools used in the manufacturing process.¹

Since these parts are often processed on highly sophisticated and accurate machine tools, subsequent removal of a fixtured workpiece from the machine and transfer to a separate inspection device seems redundant and unnecessary. A more efficient and straightforward approach would be to instrument the machine tool performing the process with adequate probes and sensors that would easily interface, both in hardware and software, to the machine and would thus utilize the inherent accuracy capabilities of the machine tool to be the verification bed for the part. Potential errors inherent in the machine to inspection device transfer (and possible repetition) such as fixturing orientation could be avoided. Obvious time savings would be realized. In fact, in a recent study by Southwest Regional Institute for a project entitled "High Productivity and Precision Machining Program" for the National Center for Manufacturing Sciences (NCMS), the highest priority focus area with 81% of the surveyed in favor was In-Process Dimensional Measurement. This survey was conducted among numerous NCMS members which include a vast majority of the most noted domestic machine tool builders.²

The successful completion of the MetreX effort will help to strengthen the U.S. machine tool industry and help to build more efficient, higher performance aircraft turbine engines at more affordable costs, and help the entire U.S. manufacturing base by making possible the economic production of complex shapes—fundamental elements in a wide range of components used for military, industrial, scientific and medical purposes.

Research Objectives

The project is directed to the development and demonstration of a Capacitive Non-Contact Analog Probe (CNAP) and a Capacitive Array Dimensional Measurement (CADM) system for inspecting complex contours. This system incorporates technology developed by Extrude Hone's MetreX Division, and will provide a fast, low cost method of measuring complex shapes. In the case of the CADM, the gauge will measure a large number of points over a surface array that is electrically scanned in milliseconds, the time required being substantially less than existing methods. The technology is being further developed to be used on the machine tool, on-site, and consideration is being given to the requirements of measurement in various processing environments. Direct software interfacing techniques

¹ Air Force Wright Aeronautical Laboratories, "AFWAL-TR-88-4177 Final Report for 1 October 1987 - 31 December 1988," *Manufacturing Technology Program Assessment of New Sensor Technologies for the U.S. Machine Tool Industry*, Sept. 1988.

² Mechanical Technology Inc., *High Productivity and Precision Machining Program: Integration Plan*, NCMS-89-PE-4.1 (Ann Arbor: National Center for Manufacturing Sciences, 1990), p. 1-3.

are being researched and if necessary, developed to provide logical and straightforward communication to new or existing machine tools.

The project is validating the ability of the capacitance sensors to accurately and repeatedly measure the complex shape of a workpiece. Mathematical modeling of the designed sensor(s) and probes is being performed and analyzed. A demonstration lab unit and specifications for the equipment design with interfacing protocol to as wide a range of machine tool controls as practical to perform the inspection procedure on the machine tool are being established. In addition, the system's ability to withstand the rigors of a hostile machine tool environment are being tested. The research will culminate in the design, fabrication and test of the entire system, the integration of the system with a machine tool process, the optimization of the process and the demonstration of the total system capabilities.

Benefits

The opportunities to improve the performance of products and to permit innovative new designs for products that are dependent on components with complex shapes are limited by the high costs of producing the components and verifying their shapes. For example, involute gears are very sensitive to gear misalignment. If the profile of the gear is manufactured incorrectly, the subsequent misalignment will cause the shift of the bearing contact toward the edge of the gear tooth surfaces and transmission errors that cause gear noise. Much time and expense is spent in the design and engineering, as well as the manufacture of these gears simply to overcome the control variations of the gear producing process. An innovation to improve the manufacturing of gears by providing timely verification and subsequent correction of the process would provide a welcome technology edge to producers, end-users, as well as U.S. manufacturers of gear machine tools.

The need to verify complex shapes and accordingly correct their manufacturing processes has been recognized by the aircraft turbine engine industry. Numerous components in a jet engine including turbine blades, impellers, and Integrally Bladed Rotors (IBR's) possess complex contours and shapes and require very detailed and complicated inspection procedures. An inspection system incorporating dimensional and tolerance data measurement in a timely fashion as close to the manufacturing process as possible would lead to reduced inspection time and costs, higher standards of accuracy, increased workpiece throughput, and reduced machining labor costs.

As devices continually and increasingly employ complex shapes and near net-processed components gain popularity, the need to verify these components' complex shapes will grow accordingly. If the proposed system of near-process verification of those shapes finds successful integration, the subsequent growth of complex shapes and their manufacturing processes would be compounded. The expanded design opportunities for new products with special capabilities will offer significant performance and strategic benefits in a range of applications including turbine engine components, gears, and cutting tools.

RELATED WORK

Capacitance sensing technology has been used since the 1950's for measuring the thickness of metal strips and coatings, the expansion and fatigue of metals, and the size, depth and cylindricity of precision bores and shafts. One of the important benefits of capacitance in these and other applications lies in its ability to measure such parameters without actually contacting the workpiece surface. Other advantages include a high frequency response, excellent linearity and resolution, and convenient portability. Typical configurations of conventional capacitance dimensional measurement systems are singular or differential in application. Extrude Hone's MetreX Division has pioneered the effort to develop capacitance technology for surface finish evaluations, and has recently developed basic capacitance sensor array capabilities for shape, edges, proximity, and slip under a Phase II Department of Energy SBIR project.

The Division is currently working on a Phase II SBIR grant sponsored by NASA to demonstrate the feasibility of integrating all of the tactile attributes into a single sensor array, to develop software to

control sensor scanning schemes and to establish adaptive grasping control algorithms.

One phase of the current Air Force project is to build on the results of this program to integrate the dimensional measurement capabilities of a capacitance array and probe to check the conformity of a complex workpiece shape to a known reference shape on the machine tool.

The Division was granted a license to patents and know-how covering the use of fringe-field capacitive technology developed at the University of Washington. The technology supplements the Division's existing techniques with surface profilometry and dimensional measurement capabilities. The Division has been working to incorporate this technology for those applications which require more than average surface roughness measurements. The advantages of current capacitance metrology including fast response time and rugged sensors for in-process inspection are enhanced and strengthened by this technology.

This technology was one of three selected for evaluation for the Automated Disk Slot Inspection System (ADSIS), an Industrial Modernization Incentives Program (IMIP) through Garrett Engine Division of Allied Signal Aerospace Company. The objective of the program was to develop an automated system of inspecting gas turbine engine disk slots using advanced inspection technologies to achieve enhanced accuracy, repeatability, and flexibility coupled with faster data acquisition. Under this program, a spherical capacitance prototype probe designed to interface to a Coordinate Measuring Machine was delivered. This single sensor probe provides noncontact dimensional data on turbine disk slots in a scanning mode without contacting the part. After initial evaluation by both Garrett and the National Institute of Standards and Technology (NIST), the capacitance probe was selected as one of the prime sensing technologies. The other technologies, conventional touch trigger and laser triangulation probes, were deemed too slow and not capable of being configured small enough to reach into the limited access areas required.

This Air Force effort is building on the results of these programs to establish the feasibility of integrating the measurement capabilities of this Capacitance Non-Contact Analog Probe to measure dimensions, features and surface characteristics in conjunction with an appropriate machine tool.

Bibliography

Garbini, J.L., L.J. Albrecht, J.E. Jorgensen and G.F. Mauer, "Surface Profilometry Based on Fringing Capacitance Measurement," *ASME Journal of Dynamic Systems, Measurement and Control*, Vol. 107(3), (September, 1985).

Garbini, J.L., F.T. O'Neill and J.E. Jorgensen, "A High Speed Fringe-Field Capacitive Hole Inspection System," *Modeling, Sensing and Control of Manufacturing Processes*, ASME, DSC-Vol.4, (1986).

Mathias, R.A. and J.L. Garbini, *An Advanced Force Sensor System for Unmanned and Flexible Manufacturing*, presented at the Sensor Technology for Untended Manufacturing Conference, Chicago, Illinois, SME Paper No. MS84-910, (1984).

Tsukamaki, T., J.L. Garbini and J.E. Jorgensen, "An Algorithm for the Control of a Bi-Stable Multi-Segmented Manipulator," presented at the Fifth IASTED International Symposium on Robotics and Automation, (November 12, 1984).

519-31

150489

P-11

N93-25580

ULTRASONIC POLISHING

Randy Gilmore
Extrude Hone Corporation
8075 Pennsylvania Avenue
Irwin, PA 15642

ABSTRACT

The ultrasonic polishing process uses high-frequency (ultrasonic) vibrations of an abradable tool which automatically conforms to the workpiece, and an abrasive slurry to finish surfaces and edges on complex, highly detailed, close tolerance cavities in materials from beryllium copper to carbide. Applications range from critical deburring of guidance system components to removing EDM recast layers from aircraft engine components to polishing molds for forming carbide cutting tool inserts or injection molding plastics. A variety of materials including tool steels, carbides and even ceramics can be successfully processed. Since the abradable tool automatically conforms to the workpiece geometry, the ultrasonic finishing method described offers a number of important benefits in finishing components with complex geometries.

INTRODUCTION

The automatic finishing of edges and surfaces of complex detailed geometries on critical cavities remains one of manufacturing's most important challenges. Finishing of molds and die cavities accounts for more than one hundred million dollars annually, primarily in hand finishing by the most skilled mold makers.

Finishing edges and surfaces of critical components for such devices as medical implants or inertial guidance gyro systems is equally time consuming. When these component geometries provide through-flow areas, abrasive flow machining can often be used—but geometries with blind cavities are not applicable or require difficult tooling. These components must be free from burrs or "smeared metal" under 20 to 50 power magnification inspection.

Uniform removal of unwanted surface layers (such as thermal recast from EDM or laser machining, layers suspected of intergranular attack, "white layers," or oxidation) selectively without masking and without losing desired detail or tolerances is another problem for manufacturers, labored over with hand work and sometimes resulting in scrap or compromised designs.

In the process of developing the ultrasonic machining process over the past several years for the purpose of machining brittle materials such as glass and graphite, it was discovered that by using a tool material that was easily abradable in ultrasonic machining—such as glass or graphite—the tool would readily shape itself to the workpiece providing a near perfect "mirror-image" conjugal form and thereby providing a uniform machining gap for ultrasonic machining. The result was that although the tool was being "machined" much faster than the workpiece (in fact, because of this) the work performed on the workpiece was uniform. Whether the objective was to polish, remove a surface layer, or deburr and lightly radius the edges, the uniformity and gentleness of the work performed retained the detail and close tolerances.

THE PROCESS

The SoneX Ultrasonic Machining System offers automatic operation, not dependent on manual polishing skills. The process uses high frequency (ultrasonic) vibrations of an abradable tool which automatically conforms to the workpiece, and an abrasive slurry to finish surfaces and edges on complex, highly detailed, close tolerance cavities in materials from beryllium copper to carbide (Figure 1). Applications range from critical deburring of guidance system components to removing EDM recast layers from aircraft engine components to polishing molds for forming carbide cutting tool inserts or injection molding plastics.

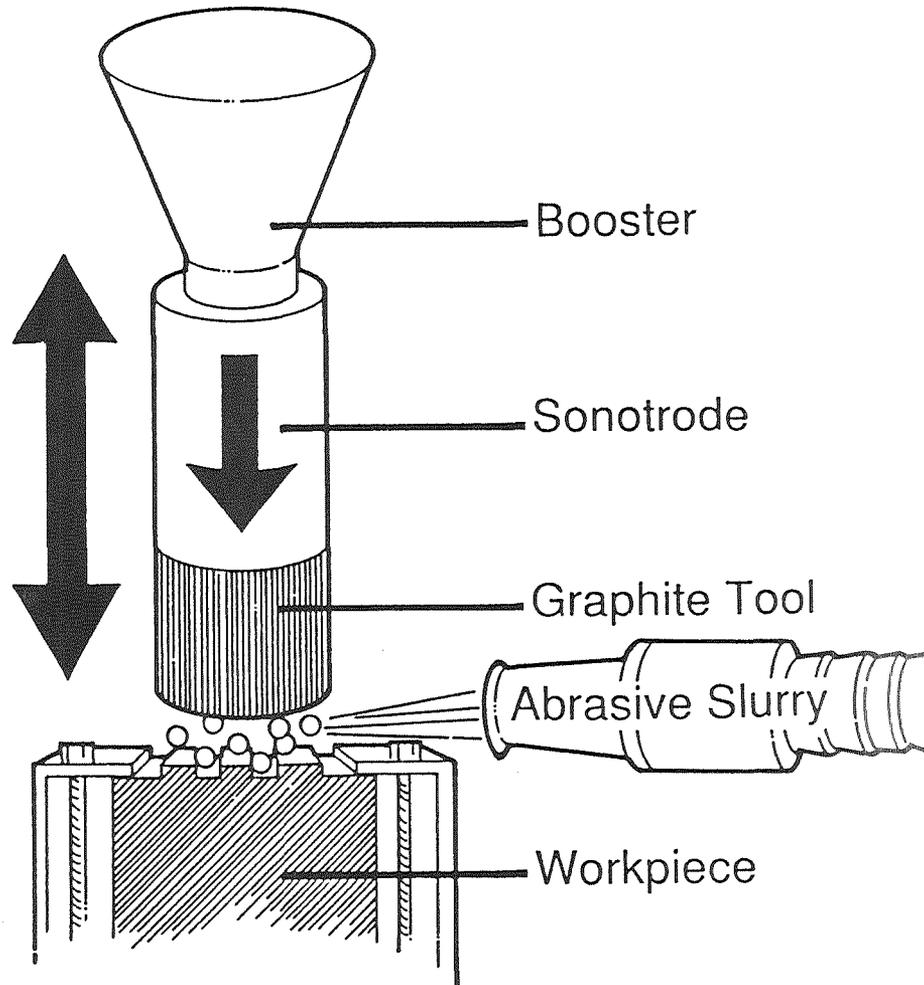


Figure 1 - Ultrasonic Polishing Schematic

The extent of polishing required is determined by the initial surface roughness of the workpiece and the finish required after polishing. Typical surface improvements range from 5:1 to 10:1; finishes as low as 4 μ inch R_a can be achieved. A variety of materials including tool steels, carbides and even ceramics can be successfully processed. Since the tool is not preshaped, but rather conforms to the workpiece configuration, indexing and registration of the tool and workpiece is not required. The part shown before and after processing (Figures 2 and 3) is a 12.7 mm (1/2 in.) diameter coining die. Ultrasonic polishing was used to remove the machining marks left by the CNC engraving operation. Cycle time is fast at under ten minutes.

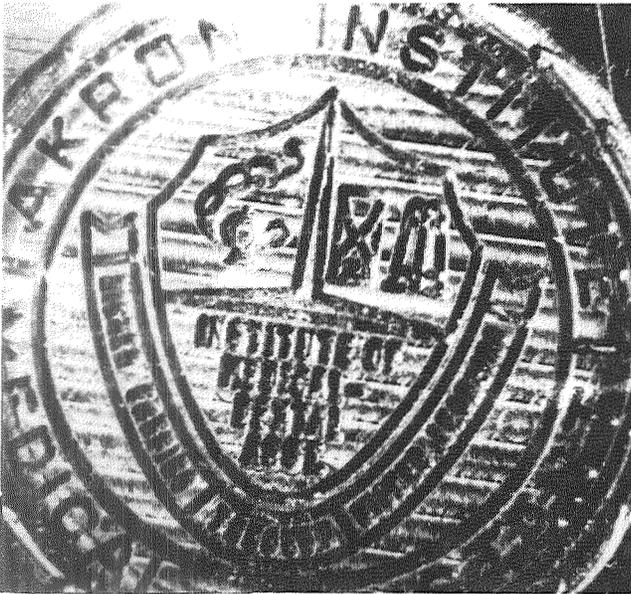


Figure 2
Coining die shown before and after processing.



Figure 3

The photomicrographs in Figures 4 and 5 compare the original ram EDM finish of 30 μm R_a on a carbide compacting with an area on the same die which has been ultrasonically polished to a final finish of 13 μm R_a , removing only 0.005 mm (0.0002 in.) of material.



Figure 4

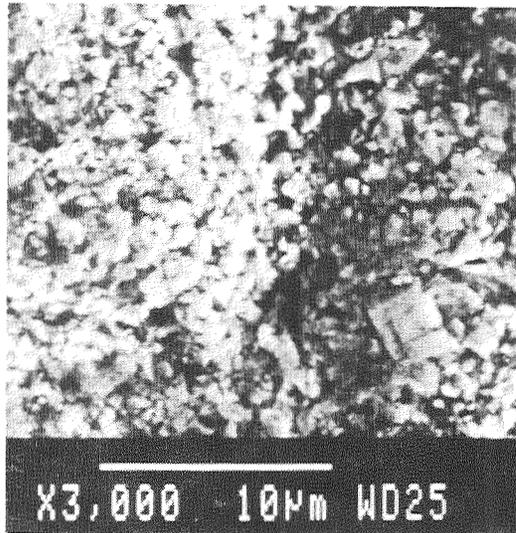


Figure 5

Photomicrographs at 3000X magnification show the surface improvement after ultrasonic polishing.

Since the abrable tool automatically conforms to the workpiece geometry, the ultrasonic finishing method described offers a number of important benefits in finishing components with complex geometries, including:

- No specially preshaped tools are required; consequently, even low volume components are applicable;
- No precision alignment of the polishing tool to the workpiece is required;

- Work is uniform across the workpiece surface.
- Surface improvement is 3:1 or more on machined, EDM'd and cast surfaces;
- Edges can be deburred and lightly radiused;
- No special operator skills are required; and
- The system operates automatically without operator involvement.

To provide the background for the research performed under the PRDA project, EHC was funded under a "Seed Grant" sponsored by the Ben Franklin Program of the Commonwealth of Pennsylvania to examine the feasibility of ultrasonic polishing. The effort yielded general ultrasonic polishing parameters which were examined more closely and optimized under this PRDA contract. The "Seed Grant" research was comprised of three tasks to optimize machining parameters, to examine various polishing and tool materials, and to evaluate the effects of ultrasonic polishing on various materials with different machined surfaces. From the data collected during this study it was clear that ultrasonic polishing offers an effective means of improving surface finish in any material of sufficient hardness. Significant results included:

- The degree of improvement is largely dependent upon the beginning surface finish; an EDM'd finish of 300 μ inch (7.5 μ m) R_a , for example, can be ultrasonically polished to a 150 μ inch (4 μ m) R_a in a relatively short cycle time—about 10 minutes. A 15- μ inch (0.4- μ m) R_a ground finish, on the other hand, can probably only be improved to about a 10- μ inch (0.3- μ m) R_a finish with cycle times approaching 15 minutes or more.
- The best abrasives include silicon carbide for aluminum and tool steels that have not been thermally machined, boron carbide for thermally machined tool steels and soft ceramics, and diamond for tungsten carbide and hard ceramics. Abrasive mesh sizes from 320 to 600 mesh yielded the best results with good machining speeds and acceptable surface finishes. Abrasive particle concentrations of 35 percent by volume are optimal for polishing. Vibration amplitudes range from 0.0004 in to 0.0015 in (0.01 to 0.038 mm) with the best frequencies achieved at the 20- to 20.5-kHz range.
- For most metals a static load of 4 to 8 pounds (1.8 to 3.6 kilograms) works best; while for ceramics and other brittle materials, static loads of 2 to 4 pounds (1 to 1.8 kilograms) achieve the best results. Good flushing is important to the success of the polishing process. The best carrier of abrasives is water.
- Simple shapes as well as complex three-dimensional shapes can be ultrasonically polished; however, the ability to successfully polish vertical sidewalls was not clearly demonstrated. The actual metal removal that occurs is typically less than 0.0005 in (0.013 mm).
- The best horn material appears to be either nonhardened tool steel or aluminum based on cost, machinability, bonding and acoustic property considerations. A cylindrical shape offers consistent vibration amplitudes and is relatively easy to manufacture. The best bonding technique uses a two-part epoxy applied after the horn was heated.
- Graphite has proved to be the best tool material offering low cost, good availability and excellent ultrasonic abrasability. A wide range of materials can be polished by ultrasonic techniques. The only limiting factor is that the material be sufficiently hard so that the abrasive particles do not impinge into the surface. The hardness of the material will affect machining speed and surface quality.

PROCESS CAPABILITIES

A complete report describing all of the activities associated with a PRDA program to establish and demonstrate prototype hardware and control software for automatic, close tolerance control of the finishing operation for complex components can be requested from WL/MTPM, Wright-Patterson Air Force Base, Ohio 45433-6533. The technical effort was conducted in one phase over a 24-month period. Objectives of the program were to optimize machining parameters and validate the ability of ultrasonic machining to provide reliable and reproducible finishing results. A prototype system was designed, fabricated and assembled as part of this effort. In addition, all aspects of operation and control were tested, the process was further optimized and the entire system demonstrated. A partnership with Kennametal and Allied Signal Aerospace was established for this program. Highlights of the program are discussed.

Deburring

For examining deburring and edge finishing parameters primary considerations included abrasive type (boron and silicon carbide), size (240 to 600 mesh) and concentration (20 to 40 percent). Samples were machined out of A2 tool steel with a 0.25-in (6.35-mm) slot milled across the face of the workpiece. Both silicon and boron carbide abrasives from 320 mesh to 600 mesh were used in tests to remove the burr. Abrasive concentration was 40 percent by volume. (In ultrasonic machining of graphite and ceramics, abrasive concentrations typically range from 12 to 25 percent.) Boron carbide performed about 50 percent faster than silicon carbide; the difference in particle size was negligible. After the full face of the vibrating tool (graphite) was in contact with the workpiece (about 10 minutes), boron carbide removed the burr in a 5-minute cycle; silicon carbide in a 12-minute cycle. The photograph presented in Figure 6 shows two of the deburring samples, one before and one after processing.

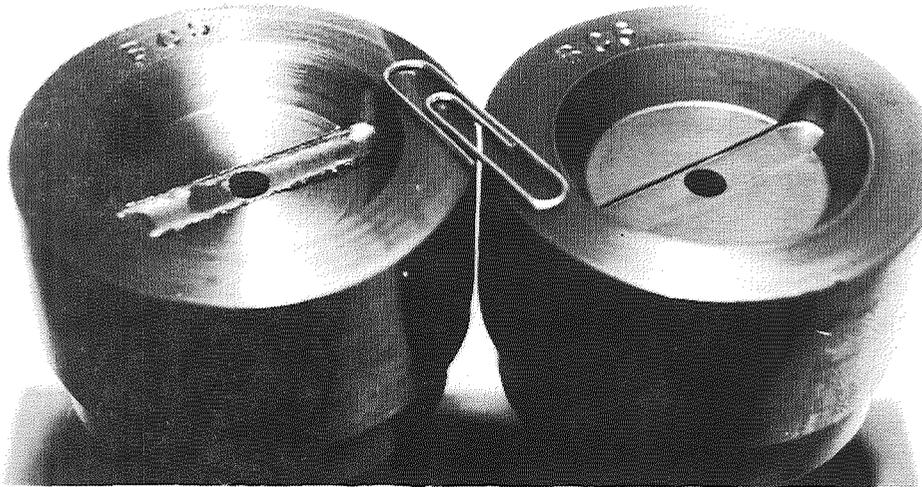


Figure 6 - Deburring Test Samples Before and After Processing

There are limitations to the burr size which can be effectively removed with ultrasonic polishing. The minimum radius produced is dependent on the burr size as well as the abrasive grain size. With 320 mesh, a 0.002- to 0.003-in (0.05- to 0.076-mm) radius was produced. For deburring, the best results are achieved with a flat piece of graphite rather than a preshaped tool.

Recast Removal

Initial recast removal tests were designed to evaluate the effects of different abrasive sizes, types and concentrations on EDM'd surfaces. Tests were performed in A2 tool steel and aluminum with silicon carbide and boron carbide of varying abrasive sizes and concentrations. The results of these tests in A2 are charted in Figures 7 and 8. Boron carbide produced finer finishes in A2 with a typical surface improvement of 60 percent compared with 40 to 45 percent for silicon carbide. Boron carbide can accomplish a greater degree of surface finish improvement per unit time, making it possible to process a higher quantity of work. In addition, boron carbide wears at a lower rate, requiring less frequent replacement. However, boron carbides costs about 10 times more than silicon carbide. In aluminum, 320-mesh silicon carbide achieved a 70 to 75 percent improvement in surface finish, most of it occurring in the first 5 minutes. Boron carbide would not be cost effective on aluminum.

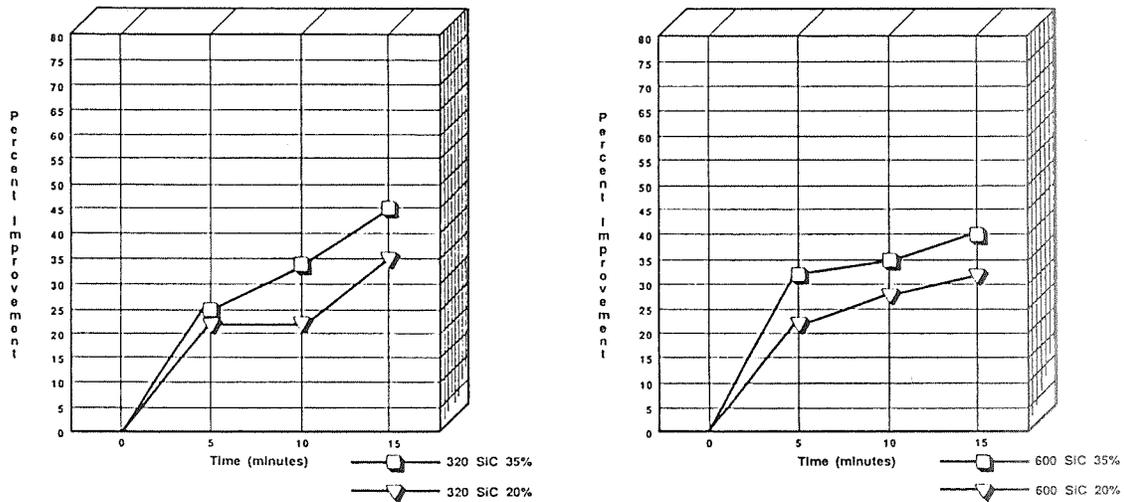


Figure 7 - EDM Recast Removal in A2 with Silicon Carbide

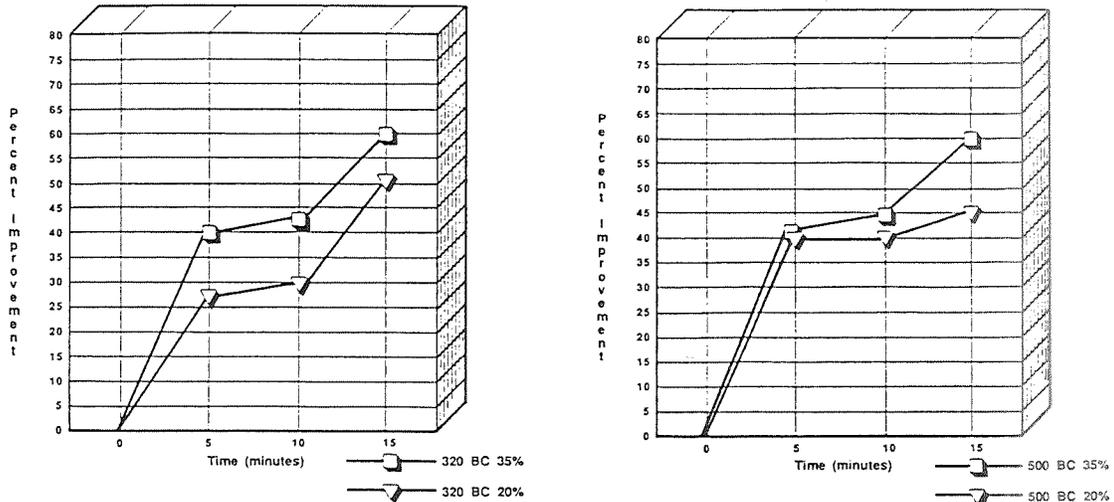


Figure 8 - EDM Recast Removal in A2 with Boron Carbide

The differences in both cutting speed and final surface improvement between the 320-mesh abrasive and the 500-mesh abrasive of boron carbide or the 600-mesh abrasive of silicon carbide are negligible. Typically 320-mesh abrasive will achieve the greatest surface finish improvement in a cycle of up to 15 minutes-after this the surface finish improvement levels off. A 600- or 1000-mesh abrasive can achieve slightly better surface finishes when cycle times are longer.

Surface finish after EDM'ing ranged from 196 to 221 μ inch (5 to 5.6 μ m) R_a in A2; 229 μ inch (5.8 μ m) in aluminum. After polishing, finishes ranged from 88 μ inch (2.2 μ m) with boron carbide to 118 μ inch (3 μ m) with silicon carbide in A2 to 60 to 70 μ inch (1.5 to 1.8 μ m) in aluminum. Typical material removal ranged from 0.0001 to 0.0003 in (0.003 to 0.008 mm) for silicon carbide to 0.0003 to 0.0005 in (0.008 to 0.013 mm) for boron carbide.

A series of polishing tests was conducted on EDM'd A2 surfaces with an average incoming surface finish of approximately 100 μ inch (2.5 μ m). Tests were limited to polishing with boron carbide abrasives because earlier tests had shown boron carbide to be superior to silicon carbide for polishing EDM'd A2 surfaces.

With 320 mesh boron carbide, considerable work was accomplished in the initial 5 minutes of polishing. After a 5-minute cycle using a 40 percent concentration of 320 boron carbide, a 45 percent surface finish improvement was achieved. As cycle times were increased, no noticeable surface finish improvement resulted. Further testing showed that 320 boron carbide at 55 percent concentrations did not perform as effectively; a 5-minute cycle yielded only a 28 percent improvement, 10 minutes resulted in a 40 percent improvement, and 15 minutes in a 45 percent improvement. These results are presented in Figure 9. Optimum abrasive concentration for EDM recast removal was established at 45 to 55 percent.

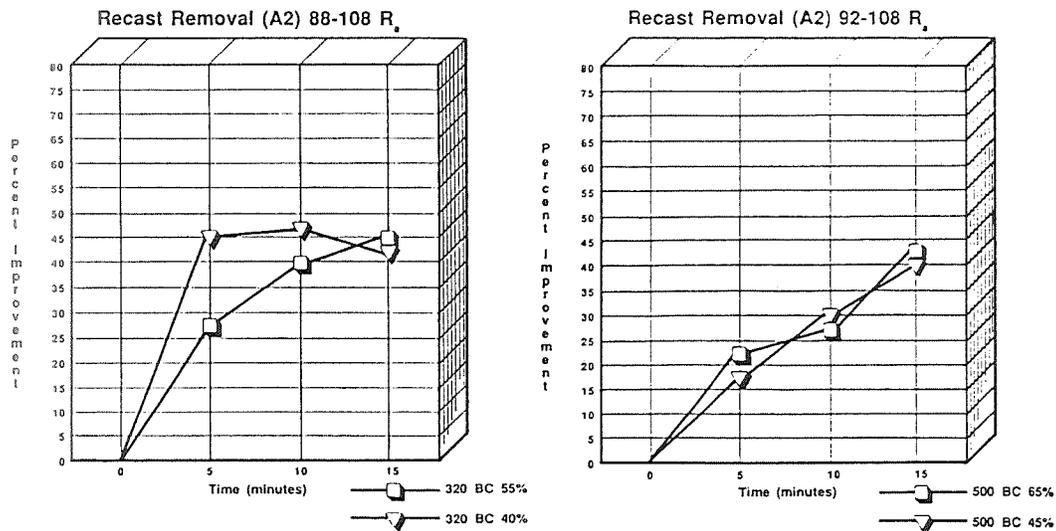


Figure 9 - EDM Recast Removal with Incoming Surface Finish of approximately 100 μ inch (2.5 μ m)

As these charts indicate, tests conducted with 500-mesh boron carbide showed an optimum surface improvement of 42 percent after a 15-minute polishing cycle with no appreciable difference in concentrations of 45 percent and 65 percent. With finer abrasive particles, longer cycle times do yield better surface finishes. Extended processing tests were conducted to establish whether longer cycle times would continue to show improvements in the surface finish. With beginning EDM'd finishes ranging from 104 to 112 μ inch (2.6 to 2.8 μ m) R_a , a 30-minute cycle yields surface finish improvements of approximately 25 percent over the initial 40 percent improvement achieved in 15 minutes. An additional 30 minutes of processing yielded only a marginal 5 percent improvement above that. While most of the surface improvement is achieved in the first 15 minutes (40 to 45 percent), additional processing of up to 15 minutes will yield a total surface improvement of 65 to 70 percent.

Residual Stress Analysis

Components that are machined with high localized energy, especially those machined by thermal processes, will commonly display stresses at the surface. These stresses are typically residual tensional stresses caused by the rapid heating and then cooling of the surface of the material being machined. In many applications, particularly for aerospace components, residual tensional stresses are not acceptable: these stresses can lead to lower cycle fatigue life because they are characterized by microcracks that when exposed to pressures or forces in use will propagate and cause cracking, which eventually could cause failure of the component.

A series of tests was conducted to examine the impact of ultrasonic polishing on residual stresses. Both heat treated and nonheat treated material were tested to determine what effect hardness had on residual stress. Both types of material were EDM'd prior to ultrasonic polishing. Several samples were then ultrasonically polished and a fewer number were left unpolished to establish the degree of stress prior to ultrasonic polishing.

As a result of ultrasonic polishing, the tensional residual stresses were negated and, in fact, changed to residual compressive stresses. This type of residual stress is advantageous to the strength and fatigue cycle of the material. Normally such compressive stresses are created by shot peening, a highly uncontrollable process. With ultrasonic polishing, residual compressive stresses can be accurately and repeatedly imparted on the workpiece. Although the depth of residual compressive stress was not as great with ultrasonic polishing as with shot peening, it is believed that this can be improved with lower frequencies and/or higher power (amplitudes). A further advantage is that while the stress characteristics are being enhanced, the surface finish of the component is being accurately improved. Lambda Research Laboratories was contracted to analyze the results of the residual stress tests.

Examine Ability to Process Larger Workpieces

True ultrasonic machining occurs at a frequency of 20,000 Hz or higher when frequencies are not audible by the human ear. Ultrasonic machining and polishing at 20,000 Hz are limited to driving a tool that has a maximum surface area of 10 square inches (64.5 square cm). To operate large tools the frequency must be lowered out of the ultrasonic range. A 10,000-Hz generator, transducer, booster and tool were incorporated into the prototype ultrasonic polishing system to permit processing of areas up to 25 square inches (129 square cm).

Primary concerns with 10-kHz polishing include soundproofing and tool tip attachment. A 10-kHz tool operating at 50 percent power can generate noise levels of 200 decibels-far above the OSHA guidelines. In addition to ear protection a soundproof room was constructed for the 10-kHz testing. Since wavelength at 10 kHz is twice that at 20 kHz, the cutting tool assembly is about twice as long as that of the assembly for 20 kHz (»23.5 in compared to 13.5 in {597 mm compared to 343 mm}). This added length and larger tool make it more difficult to achieve a good bond between the tool and tip. Attaching tool tips for 20-kHz polishing is successful 80 to 90 percent of the time, while 10-kHz tool tips can be attached properly only 30 to 40 percent of the time.

Sonic (10 kHz) polishing uses the same amplitude as ultrasonic (20 kHz) polishing, but since the surface area is considerably larger in sonic polishing, more work can be accomplished in unit time. Results from the testing performed in this program show that areas up to 25 square inches (129 square cm) can be polished simultaneously. This area can be one large workpiece or a series of smaller workpieces. Testing also indicates that the resultant surface finish is comparable to that of 20-kHz polishing.

Although the feasibility of polishing larger areas at 10 kHz has been demonstrated, further testing is required to optimize polishing at this frequency. Future investigation should be aimed at the effects of 10-

kHz sound on the human ear and how to minimize any possible hazards to the operator. Additionally, a more reliable bonding mechanism needs to be established to facilitate 10-kHz processing.

Prototype Ultrasonic Polishing System

The prototype ultrasonic finishing equipment was designed in four subsystems including mechanical components, drive and controller systems, ultrasonic hardware and slurry system. The basic design of the superstructure of the machine is a four-post configuration for maintaining slide accuracies of 0.0001 in (0.003 mm). The XY worktable is capable of holding workpieces up to 18 x 24 x 12 in (457 x 610 x 305 mm) and weighing up to 350 pounds (160 kilograms). All three axes (X, Y and Z) are driven by AC servo motors with positioning capabilities of ± 0.00025 in (0.006 mm) through the full travel envelope. The Kurt Robocon II controller was chosen for integration due to its ability to accept analog voltage input commands from a series of force transducers built into the tooling head; the ability to use a PC allowing customized screen generation and NC G-code programming; and the ability to control up to six axes in synchronization.

The ultrasonic generator is capable of automatic resonance search and following allowing constant scanning of the transducer/tool assembly to continuously adjust and optimize frequency. In addition, the output power level can be adjusted from 0 to 1 kw to permit different amounts of stroke (or amplitude) to be used. Finally, the generator can be operated in both 20 and 10 kHz frequencies for processing larger workpieces.

The slurry system can deliver up to 10 gallons (38 liters) per minute to the machining area with a tubing pump incorporated for low wear and easy maintenance. The holding tank is tapered to prevent abrasive packing and a chiller is used to maintain slurry temperature.

Selected parameter data were incorporated into a menu-driven display with prepackaged programs for automatic selection of machining parameters based on the depth of the area to be polished, beginning surface finish and previous machining method. Performance requirements for the workpiece, tool, machine and controller were specified.

Based on the design specifications and performance requirements, the control and machine tool were built and integrated. The control system is comprised of an IBM compatible industrial PC, monitor and keyboard; the Kurt Robocon II controller with a special analog input board; motion control amplifiers and servos; a 16-position auxiliary input/output board; Heidenhain linear encoders with times five multiplier boxes; a SLICE 10/20 kHz, 1 kw ultrasonic generator with special interface board; operator controls and pendant workstation; and transformers for the various components. Special consideration was taken to ensure easy accessibility for trouble-shooting and maintenance.

The die set, frame assembly and weldments were subcontracted to outside vendors. The XY table was purchased from Setco and incorporated special features allowing accurate mounting to the machine body. All cover and guard assemblies were fabricated at EHC and all assembly was also performed at EHC. As with the control system, special consideration was given to accessibility and ease of maintenance. All cover assemblies were manufactured from Alucobond panels that can easily be removed and replaced by a single individual. Assembly of the machine tool required approximately 14 weeks after receipt of all components and ran concurrently with the completion phase of the control system.

After completion of the control system and machine tool, integration of these subsystems as well as the slurry system was begun. During integration, software control over input/output functions was verified, push button control over input/output functions was tested and slurry system functions were completed. Total time for the integration phase was approximately 10 to 12 weeks. The prototype system is pictured in Figure 10.



Figure 10 - Prototype Ultrasonic Polishing System

The parameters selected for testing of the completed prototype ultrasonic finishing machine included mechanical accuracy, operating software and operator interface, control hardware and polishing performance. The XY slide was found to position within ± 0.0002 in per foot (0.005 mm per 305 mm), tracked straight within ± 0.0002 in (0.005 mm) and was repeatable within ± 0.0001 in (0.003 mm). The moveable platen positioned within ± 0.0002 in (0.005 mm) of commanded position and repeated within the same tolerance. Parallelism to the worktable surface was not acceptable at the testing phase, so corrective action was implemented. The operating software was tested for the responsiveness of the servo loop interaction with load cell feedback. At the time of testing it was determined that the servo loop was not tight enough and corrective action was recommended and implemented. Control hardware wiring, particularly voltage levels and grounding, was verified. Safety issues were examined and accepted.

Testing of polishing parameters revealed that the servo loop problem impacted the stability of machining. It was decided that the corrective actions recommended for the servo loop would alleviate the problem of stability. Polishing speed was examined and rates as high as 0.030 in (0.762 mm) of vertical travel per minute were realized. Polishing quality testing showed surface finish improvements as much as 7:1, exceeding the expected results. Parallelism of the moveable platen was determined to be unacceptable due to poor timing between the two ballscrews. The timing was reset and parallelism was improved to ± 0.0002 in (0.005 mm) over the full platen surface. The servo loop/stability of machining problem was found to be caused by the algorithm that controls this loop. A new algorithm was written to correct this problem.

Process repeatability was examined based on improvement of surface finish and amount of stock removal. Stock removal ranged from 0.0001 to 0.0005 in (0.003 to 0.013 mm), dependent on incoming surface roughness. Stock can be removed accurately within 10 percent of the stock removed; stock can also be removed to a desired depth within 20 percent and is repeatable from workpiece to workpiece within 10 percent. Limitations observed are confined to vertical or near-vertical side walls where stock removal was as much as 50 percent less than that of frontal surfaces.

Surface finish improvements are the greatest when incoming surface roughness is high. In this case, surface finish improvements as high as 10:1 were accomplished. On components with beginning surface finishes in the range of 100 μ inch (2.5 μ m) R_a , improvement averaged 5:1. Repeatability of surface finish improvement measured \pm 10 percent.

Process parameters were optimized including slurry, feed/speeds, polishing tip material, machining pressures and ultrasonic generator settings. Boron carbide is the best multi-purpose abrasive for ultrasonic polishing. With a 600-mesh abrasive, surface finishes as low as 8 to 12 μ inch (0.2 to 0.3 μ m) R_a are possible; finishes as low as 14 to 18 μ inch (0.36 to 0.46 μ m) R_a can be produced using 320-mesh abrasives. Optimum abrasive concentration for most applications occurs between 18 and 22 percent by weight. Polishing tips of graphite show the most promise because they rapidly conform to the workpiece configuration and are of relatively low cost. Graphites with particles sizes of 0.00004 to 0.0001 in (0.001 to 0.003 mm) performed the best. The optimum static pressure ranges between 1 and 1.5 pounds per square inch (0.5 and 0.7 kilograms per 6.452 square centimeters).

The prototype system has been demonstrated at three international machine tool shows and has been discussed in numerous technical presentations and articles. The first commercial installation has been effected for polishing ram tips used in the compacting of tungsten carbide cutting inserts. In addition a variety of contract ultrasonic finishing is being performed at EHC.

SUMMARY

The most labor intensive, uncontrollable area of production remaining in the manufacture of precision parts involves the final finish machining operations, which frequently absorb more labor time than machining. Proper finishing of edges and surfaces affects more than the appearance or feel of a product; controlled, consistent edge and surface finishing can dramatically improve product performance and life while reducing direct labor costs. These finishing operations have been identified as the single greatest hurdle remaining in fully automating the production of precision components. By applying ultrasonic machining techniques, a process has been established for automatable, repeatable, uniform polishing with no chemical or electrical alterations in the surface. The process is applicable to a wide range of materials, including ceramics, composites, new alloys and plastics, with a level of accuracy not previously achievable. In addition, the process has the capability for integration with a tool changer for automatic loading of the workpiece and tool for uninterrupted cycling and incorporation into a flexible manufacturing cell, providing consistent, uniform and repeatable results and yielding improved component performance in this final and critical phase of the complete manufacturing cycle.

omit

**MICROELECTRONICS/OPTOELECTRONICS
PART 1**



TWO- AND THREE-DIMENSIONAL HIGH PERFORMANCE, PATTERNED OVERLAY
MULTI-CHIP MODULE TECHNOLOGY

Capt James Lyke
Wafer Scale Electronics Engineer
Phillips Laboratory (USAF)
3550 Aberdeen Ave SE
Kirtland AFB, NM 87117-5776

520-33
150490
P 10

ABSTRACT

A two- and three-dimensional multi-chip module technology has been developed in response to the continuum in demand for increased performance in electronic systems, as well as the desire to reduce the size, weight, and power of space systems. Though developed to satisfy the needs of military programs, such as the Strategic Defense Initiative Organization, the technology, referred to as High Density Interconnect, can also be advantageously exploited for a wide variety of commercial applications, ranging from computer workstations to instrumentation and microwave telecommunications. The robustness of the technology, as well as its high performance, make this generality in application possible. More encouraging is the possibility of this technology for achieving low cost through high volume usage.

INTRODUCTION

Size, weight, power. To a space system, they represent cost and complication. Reducing them, at least as far as the electronic subsystems are concerned, became the central focus of the Phillips Laboratory (PL) Wafer Scale Integration (WSI) program, sponsored largely by the Strategic Defense Initiative Organization (SDIO). At the same time, the Defense Advanced Research Preparedness Agency (DARPA) recognized the tremendous potential of such a technology to improve the performance of electronics systems, military and commercial. When DARPA and PL initiated programs in the mid-1980's to explore these new possibilities, they participated in developing a new approach within an emerging class of technologies called *multi-chip modules* (MCMs), which themselves are a part of a broader class of technologies known as *wafer scale integration* (WSI) technologies. The MCM technology discussed in this paper, known as *High Density Interconnect* (HDI), developed by complementary support from DARPA and PL with the General Electric company, is a novel technique for re-assembling bare integrated circuits in a manner that dramatically improves size, weight, and performance [1]. Perhaps more novel is the potential of this technology to address an exceptionally large class of applications: military and civilian, space and terrestrial, strategic and commercial.

BACKGROUND

Traditionally, an integrated circuit (IC) chip is placed in a single chip package (SCP) to protect it from mechanical damage and to provide electrical access to its tiny electric terminals. The problem associated with conventional electronics assembly methods based around the SCP is that they exact significant penalties upon the size, weight, power, performance, and reliability of an electronics system. Furthermore, the reliability of a system built in a manner is also non-optimal, due to extraneous materials and structural interfaces, each representing an additional failure opportunity and thermal barrier. Finally, SCPs limit the electrical design complexity that is projected for systems near and beyond the year 2000. It is suggested that for these systems, many hundreds if not thousands of signals will emanate from the terminals of individual IC chips [2].

Phillips Laboratory, in the mid-1980's, launched the Wafer Scale Integration (WSI) program in search of better packaging approaches. Here, a special emphasis was placed on technologies that were suitable for

space strategic missions. Some of the missions envisioned for the SDIO were based on capabilities that were prototyped with large brassboards. For these missions to be practical, the same functionality demonstrated in the brassboards, sometimes one or more racks full of electronics, had to be compressed into compartments that were measured in cubic centimeters rather than cubic feet.

MCM approaches can be delineated into patterned-substrate and patterned overlay approaches (Figure 2). In the *patterned-substrate* approach (Figure 2a), the interconnections between components are formed *a priori*, like a micro-printed circuit board. In the *patterned overlay* approach (Figure 2b), the approach used for the HDI process, the interconnections are not formed until all components are placed within the substrate.

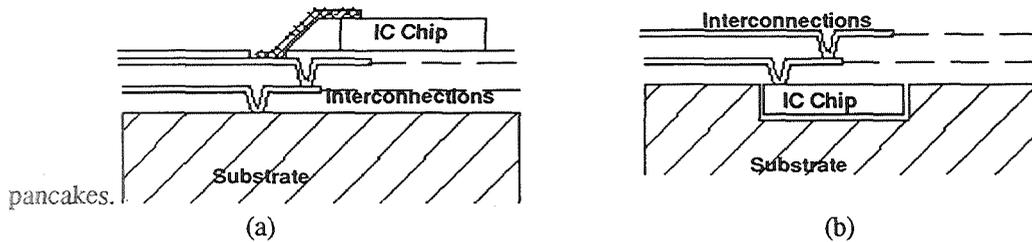


Figure 1. Hybrid wafer scale integration/multi-chip module approaches. (a) Patterned substrate. (b) Patterned overlay.

THE HIGH DENSITY INTERCONNECT (HDI) PROCESS

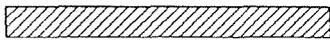
A simplified version of the sequence used to create HDI modules is depicted in Figure 2. Construction of an HDI module starts with a blank piece of substrate material of suitable flatness and quality. The substrate is processed first by forming wells or recesses into which the components are placed. Different component thicknesses are accommodated by recesses of different depths, formed in such a way as to yield an assembly that is essentially planar (within 0.0005 inches) after all components are placed. Components are mounted into the substrate with computer-controlled component placement equipment that uses the opposite corners of each component as alignment fiducials. Components are secured to the substrate with a special thermoplastic, chosen for its thermal conduction, thermal coefficient of expansion, and flow viscosity properties. Aluminum can be sputtered on areas of the substrate not already covered by components to facilitate the longer direct current interconnections (such as power and ground) and the terminals of the HDI modules.

Subsequent steps of the HDI processing sequence relate to the novelty of the formation of the so-called *patterned overlay*. This sequence begins with the lamination of Kapton dielectric to the entire substrate. Tiny opening holes or *vias* to the surface below the lamination are laser-formed and metallized to complete electrical connection to the ICs below. The metal system consists of sputtered titanium and copper, electroplated copper, and sputtered titanium. To pattern the metallization configuration, a sprayed photoresist layer is laser-exposed as required with an adaptive lithography system. The adaptive process is used to dynamically correct for the slight but inevitable errors that occur in component placement. Eventually, 2 to 4 additional layers of this dielectric-metal system are formed onto the HDI assembly, depending on the wiring capacity demands of a particular application [3].

Final packaging of HDI modules for military and space applications typically involves using a laser- or seam-welded kovar flat package to hermetically encapsulate the HDI substrate. These packages interface electrically to the HDI substrate through lead wires, which pass through insulating glass beads in the wall of the package. Simple wirebonds between the HDI substrate and the leads facilitate the electrical connections required. Other hermetic and non-hermetic packaging methods have been explored and remain an active area of research. These include: (1) the *integral package*, which uses the substrate itself as a final, hermetic package, eliminating *all* wire-bonds; (2) non-hermetic, direct attach techniques,

which interface inverted HDI modules directly to a PWB using an appropriate interposer (e.g., cinch connector, elastomeric, etc.); and (3) special high pin count package designs, extendible to three-dimensional HDI modules [4].

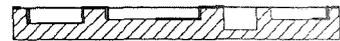
PREPARING SUBSTRATES



The HDI process begins with a flat blank of a starting material, such as alumina (most commonly used), aluminum nitride, silicon, glass, etc. This flat blank becomes the substrate, which provides a mechanical supporting structure for the HDI module.



Pockets are formed in the substrate using industrial computer-controlled milling equipment (other high volume production techniques can be implemented). These pockets become receptacles for the various integrated circuits and passive components required for the functional HDI module.



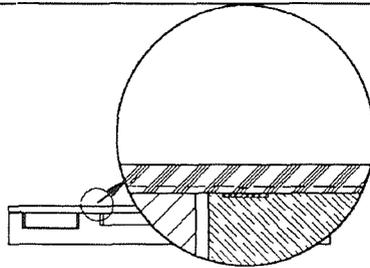
A thin layer of aluminum is deposited uniformly onto the substrate. The metal is then selectively patterned and etched to form "backside metal" contacts for certain components and other special functions.

PLACING COMPONENTS

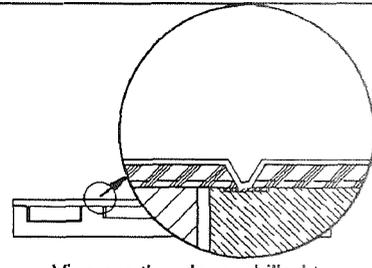


After a computer-determined quantity of adhesive is automatically distributed in each pocket, the components are transferred from dispensers with a robotic "pick-and-place" machine. The chips are placed with their electrical contact pads facing upward.

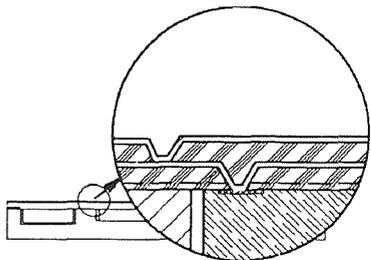
FORMING A PATTERNED OVERLAY



Following a spray-on application of Ultem 1000 thermoplastic dielectric, the first Kapton layer is laminated to the substrate at 300° C.

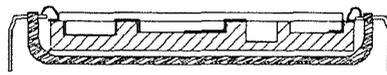


Vias are then laser-drilled to open contacts to the component terminals. A 4 micron thick metal system (Ti-Cu-Ti) is processed through sputtering and electroplating, forming the first layer interconnections.

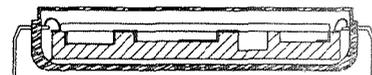


Additional dielectric layers and metallization layers are formed as necessary. These subsequent laminations utilize a thermosetting siloxane-polyimide adhesive.

FINAL PACKAGING



One of the most common modes of packaging for HDI modules is the hermetic package. In this case, the HDI substrate is glued into a package made of Kovar. Wire bonds are formed between substrate and package leads.



Finally, a package lid is welded or soldered to form a hermetic seal. Other HDI packaging options include hermetic integral substrate and non-hermetic carriers.

Figure 2. The two-dimensional HDI process.

Process Repair and Component Pre-Test

The ability to repair any MCM process is important for two reasons: (1) the value of the collection of components committed to an MCM is sometimes very significant, particularly for military applications; and (2) component yields are often too low to produce a first-pass functional MCM. Repair of the HDI process is routinely accomplished by removing the overlay, selectively removing and replacing bad components, and rebuilding the overlay. Recent tests have demonstrated as many as eleven consecutive repair cycles on individual HDI modules with no measurable differences in the physical or electrical performance characteristics of the components. Several "known good die" techniques can be employed with the HDI process to mitigate the need for module repair. One technique involves the burn-in on partially completed HDI assemblies. Unlike patterned substrate approaches, HDI allows interim testing to be performed after the metallization of each layer, if desired. A simpler technique for high volume production involves the use of a dedicated electronic membrane, built from HDI itself and configured to interface with all critical substrate components upon (temporary) contact to the substrate. Finally, a more elegant technique that mimics the conventional IC burn-in approaches can be used, which involves the use of a special adaption of the HDI process to actually create recoverable ICs. This *temporary interconnect* process, forms single-layer temporary patterned overlays on individual ICs, which can be removed after the IC is burned in. Each of these testing techniques represents an interim solution for a problem that will ultimately be solved by the semiconductor vendor.

Three Dimensional HDI

Recent SDIO-sponsored research under the PL program has resulted in the development of a three-dimensional extension of the HDI process. This "3D HDI" is unique among the various approaches that have been proposed to achieve a three-dimensional packaging system in that it is based on direct extensions of the two-dimensional HDI process. A simplified sequence for achieving 3D-HDI is shown in Figure 3. 3D-HDI combines a collection of identically-sized two-dimensional HDI modules into a very compact assembly through direct stacking. In this case, electrical contacts are formed in on the edges of individual HDI modules to be combined into a three-dimensional assembly. After the stack of HDI modules is laminated together, new HDI patterned overlays are created, which interconnect the edges of individual layers together, like a miniature backplane. Depending on the module level wiring capacity requirements, two or all four edges are utilized for patterned overlay interconnection. Several demonstration three-dimensional HDI modules have been constructed, including interconnectivity modules (Figure 4a), thermal profiling modules, and functional memory modules (Figure 4b-4d). Final packaging can be accomplished using kovar packages similar to those previously described for 2D-HDI.

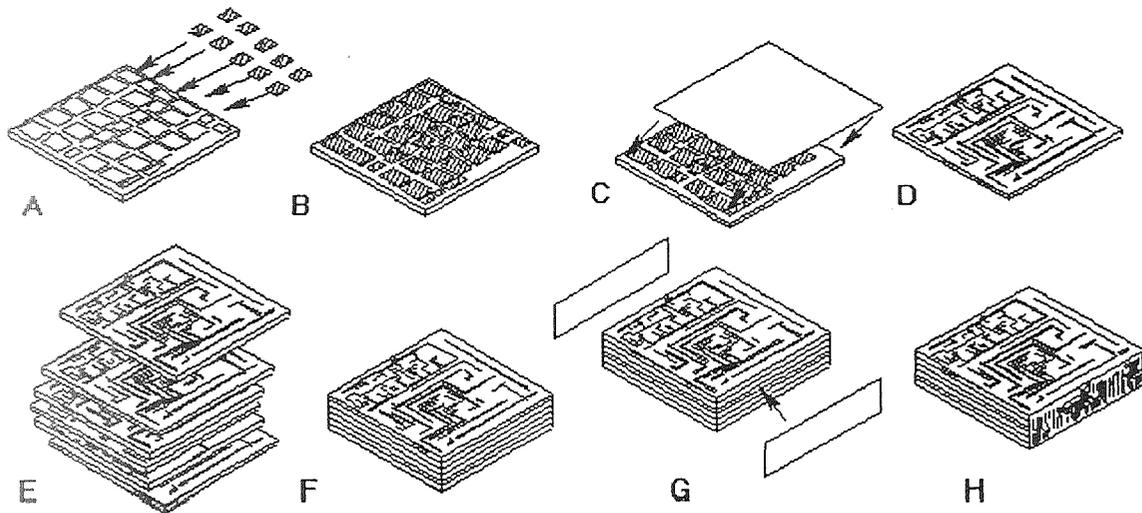


Figure 3. Three-dimensional HDI process.

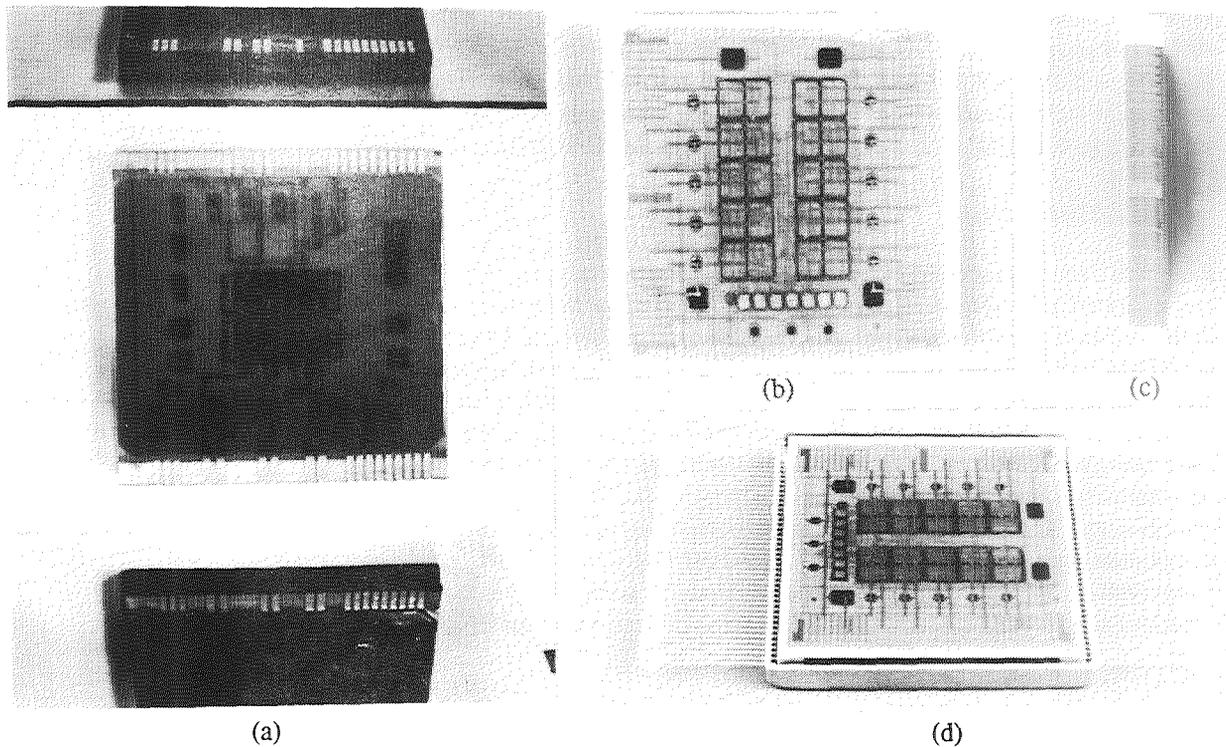


Figure 4. Three-dimensional HDI modules. (a) Interconnectivity Module. (b) Memory Module. (c) Edge view of memory module. (d) Memory module in final package.

Properties of HDI

This section addresses further the electrical and mechanical properties of the HDI process that make it viable for general application and discusses reliability and manufacturability issues.

Electrical Performance

Perhaps the greatest motivation for the development of MCM technologies is electrical performance. As the performance of digital systems built with silicon ICs exceeds 100 MHz, high fidelity interconnections structures are required. The interconnections of a PWB are reasonably high in performance, but their physical span, resulting from the use of SCPs for a conventional system forces the designer to implement transmission line design rules at a sufficiently high frequency. Although transmission line behavior is well understood, the prospect of re-designing hundreds of digital signal lines to accommodate controlled impedances is formidable. Converting such designs to MCM technologies such as HDI allows the system designer to create better systems while still employing the simpler, lumped circuit analysis and design procedures. With HDI in particular, the quality of the interconnection system is very high. The thick copper interconnections minimize series loss, and the relatively thick, low permittivity, planar intermetal dielectrics provide a relatively low capacitance. In this respect, HDI interconnections are superior in quality to the lossy interconnections that are characteristic of the integrated circuit itself.

As promising as HDI is for digital circuits operating in the lumped element regime, it is even a more enabling solution for digital and analog circuits operating in the transmission line regime and even beyond. Measurements recently performed by the Mayo Foundation indicate the capability of two-inch

HDI substrates to support impressed digital signals well above one gigahertz [5]. One of the most significant advantages in electrical performance for HDI is that the high quality interconnection structure terminates directly to the IC itself. The physically smooth transition of an HDI via is particularly important at microwave frequencies, where low-inductance, high quality terminations are desired. As such, HDI provides much higher quality transitions than those provided by wire bonds, tape automated bonding, and even solder (as used in controlled collapse chip interconnection approaches). To better understand these advantages, an effort was recently undertaken by GE to compare ordinary microwave structures built in HDI to those built with more traditional technologies.

Mechanical Design

Patterned overlay processes offer several significant advantages. First, the thermal and electrical paths are separated, providing a more efficient utilization of the available area. All contacts to the components in the substrate are formed simultaneously during assembly, using high yield fabrication methods developed originally for the construction of ICs. The patterned overlay design also supports maximal component densities. In some cases, more than 90% of a substrate area is covered by components, compared to the 6-20% in conventional packaging approaches. The interconnection capability of a patterned overlay process is great compared to MCM approaches based upon wire bonds or tape automated bonding, which are predominantly limited to connecting only the peripheral regions of components. The ability to support interconnection of contacts distributed throughout the surface of a component is promising for future generations of ICs, which will require greater communication bandwidths. The patterned overlay is more compact vertically, due to its planar construction. This planarity allows very rugged and compact three-dimensional structures to be formed.

Reliability

The robustness of HDI has been graphically demonstrated in a series of tests by PL and GE to test the performance of HDI in extreme thermal, mechanical, and nuclear environments. A summary of these tests are presented in Table 1. Central to the design of the HDI process was the selection of appropriate materials, since many failure mechanisms are related to the interactions that can occur at the interfaces of dissimilar materials. The compatibility of these materials are often exercised by the batteries of tests prescribed in military standards and specifications. For example, several thermal shocks between two temperature extremes tests the differences in thermal expansions of the various materials used in the HDI process. Underground nuclear tests demonstrate the ruggedness of a technology to withstand high amounts of total ionizing dose and thermomechanical shock. On these and other bases, the HDI process has been demonstrated quite spectacularly.

The research to further refine the understanding of HDI reliability continues. One effort, a joint venture of reliability analysts within the Department of Defense and NASA, is exploring the value of the so-called traditional military reliability tests to adequately exploit reliability problems in various MCM technologies, including HDI. This "Reliability Technology" or "RelTech" effort is conducted through theoretical modeling, test, and destructive and non-destructive analysis of groups of HDI test structures of representative complexity. These structures are instrumented with reliability monitoring structures (miniature test ICs developed by Sandia National Laboratories), as well as representative digital logic, memory, analogy, and transmission line circuits. The findings from this effort may lead to future refinement of the criteria used by NASA and the military to evaluate, qualify, and certify MCM technologies for military applications and those commercial applications which rely on military specifications for guidance in technology selection (e.g., heart pacemakers, etc.).

Manufacturability

For both engineering prototypes and high volume production, HDI has several significant manufacturability features. For low-volume production, flexibility is important. In some MCM approaches, significant non-recurring expenses (NRE) are required for the formation of many

photolithographic masks and often the construction of specialized tape automated bonding (TAB) frames. For low volume prototype fabrication, the HDI process, which is maskless and requires no TAB frames, can be particularly cost effective. Engineering changes can and have been performed overnight, in contrast to other approaches, such as cofired ceramics, which often require substantial turn-around times for even minor patterning changes. At the high-volume end of the production spectrum, the impact of NRE is, of course, amortized across a much larger number of units. While primarily a maskless process, HDI can also employ mask-based processing, due to newer, high-accuracy automated die placement equipment. With other production enhancements and the advantage of the patterned overlay process, which forms all interconnections to components as a by-product of fabrication (in contrast to patterned substrate approaches), HDI can potentially realize one of the lowest per-unit costs of any MCM process produced in volume.

Table 1. Representative test data on two-dimensional High Density Interconnect (HDI) substrates (after [1]).

Category	Test	Description
Mechanical	Constant acceleration	7,000 G
	Drop Shock	1,500 G, 0.5 ms
	Extreme mechanical shock	68,000 g linear 178,000 g centrifugal
Thermal	Thermal Cycle	1,000 Cycles -55 ^o to 125 ^o C
	Severe Acceleration	1,000 Cycles -200 ^o to 155 ^o C
Radiation	Thermomechanical Shock	Underground Nuclear Test
	Total Dose	Tested above ground to 40 Megarads total dose

APPLICATIONS OF THE HDI TECHNOLOGY

Since the inception of the PL WSI program, the interest in MCM technologies has become widespread, both in the government and private sector. Clearly, an increased level of integration allows higher performance systems to be constructed, some of which would be impossible to assemble with conventional packaging approaches. The two key considerations for candidate applications are those that would most dramatically benefit from: (1) a overall consistent reduction in interconnection path lengths in a highly controlled manner and (2) a greatly improved wiring capacity. In other cases, 2D- and 3D-HDI can facilitate the assembly of an existing system in a form compact enough to make it useful in ways not previously conceived. Finally, in those applications where "more is better", 2D- and 3D-HDI allow systems of much greater capacity in an equivalent weight and volume. A sampling of those applications that would most dramatically benefit from these traits are presented in the following paragraphs.

Magnetic Recording System Replacements

Every space system, aircraft, and many ground systems require recording systems for data storage. For space systems, an orbit's worth of data must be collected for relay to a ground station located in one position on the surface of the earth. Flight recorders are required in all military and commercial aircraft to hold key information for crash investigations or routine reliability surveys. Finally, the use of data recording systems for instrumentation is widespread in many systems. The traditional use of magnetic recorders becomes less attractive, because of the limited reliability, data recording speed, and ease of

access. For these reasons, solid state storage is being actively considered by the military and NASA. HDI has already been demonstrated in several different memory module designs, fabricated for an upcoming satellite experiment (Figure 5a). Additionally, ruggedized memory module of two different designs have been fabricated for the Defense Nuclear Agency. HDI, particular the 3D-HDI, gives semiconductor memory storage a much greater density, making the solid state storage alternative a lucrative one.

Medical Imaging System Enhancement

A medical imaging system is only one example of a system where the ability to rapidly acquire and process data is an important consideration. One of the key components of any data acquisition system is the analog-to-digital convertor (ADC), which samples a continuously variable signal and provides it a computer in a digital form. It is possible to construct fast, coarse ADCs or higher resolution but slower ADCs with existing monolithic IC processes. While in principle it is possible to assemble fast ADCs with high resolution from a group of low resolution ADCs, the variability in electrical path lengths create electrical skew and differential resistance between channels render this approach impractical when using conventional packaging methods. The first 8-bit video rate flash converters were not routinely used, for example, until companies such as TRW were able to build them monolithically. With HDI, however, the construction of multi-component ADCs may now be practical, since electrical path lengths and electrical contact terminations are directly and precisely controlled through layout. This feature, along with the low permittivity dielectrics used in the patterned overlay, create the electrical performance characteristics needed to build composite, high-performance ADCs.

Mixed-signal applications

More generally, the construction of superior mixed signal (analog plus digital) systems, such as dc-to-dc power convertors are made possible through HDI. While excellent monolithic, mixed-signal IC processes exist today, it is typically not possible to create both analog and digital devices of superior quality and density within the same IC processes. HDI allows the mixture of a wide variety of device technologies. In the same module, optimal digital (e.g., sub-micron CMOS), analog (e.g., high quality bipolar), and microwave (e.g., GaAs) IC processing technologies can be present. The superior interconnection design afforded by the patterned overlay interconnection system allows the construction of functional blocks that have a performance superior to that in a monolithic or conventionally packaged multi-component approach. The performance advantages of this capability can be considerable. For example, GE, in their own internal research, have developed dc-dc power convertors that approach 100 W/in³ density with an efficiency above 85%. Monolithic dc-dc power convertors do not presently achieve this efficiency because of the slight compromises made in the mixed-signal monolithic IC processes which accumulate to an overall lower efficiency at the system level.

Telecommunications and Microwave

The increasing emphasis on high-speed telecommunications is driving the performance of electronics. One example of a high-speed telecommunications product is the radio-frequency (rf) modem, which sends information normally transmitted through telephonic or direct-line media through radio transmission. Presently, rf modems are bulky and expensive, and will be prevalent only when they are made compact and inexpensive. They represent one of many opportunities in the telecommunications field where a high-performance MCM technology such as HDI can accelerate use through size, weight, and even cost reduction. The extension of HDI to other applications in the microwave region is also promising, due to its superior electrical performance.

Computers

Faster computers can be constructed using HDI. In several early programs with DARPA, HDI was used to construct a computer based on Texas Instruments TMS320C25 16 bit fixed point microprocessors. The HDI version of a computer previously built on PWB was shown to achieve in some cases above 80 MHz performance, substantially higher than the IC's rated performance of 40 MHz. Other experiments at PL demonstrated a measured performance of up to 71 MHz from memory components rated as 40 MHz components. More recent experiments with Lawrence Livermore National Laboratory on R3000 based computer cores (Figure 5b) built in HDI demonstrated at least a 70% performance improvement. In this case, several of the 16 MHz units (54 components, including an R3000 central processor, R3010 floating point unit, and 256 kilobytes external cache) were functionally operating up to the 27.5 MHz limit of the testing configuration.

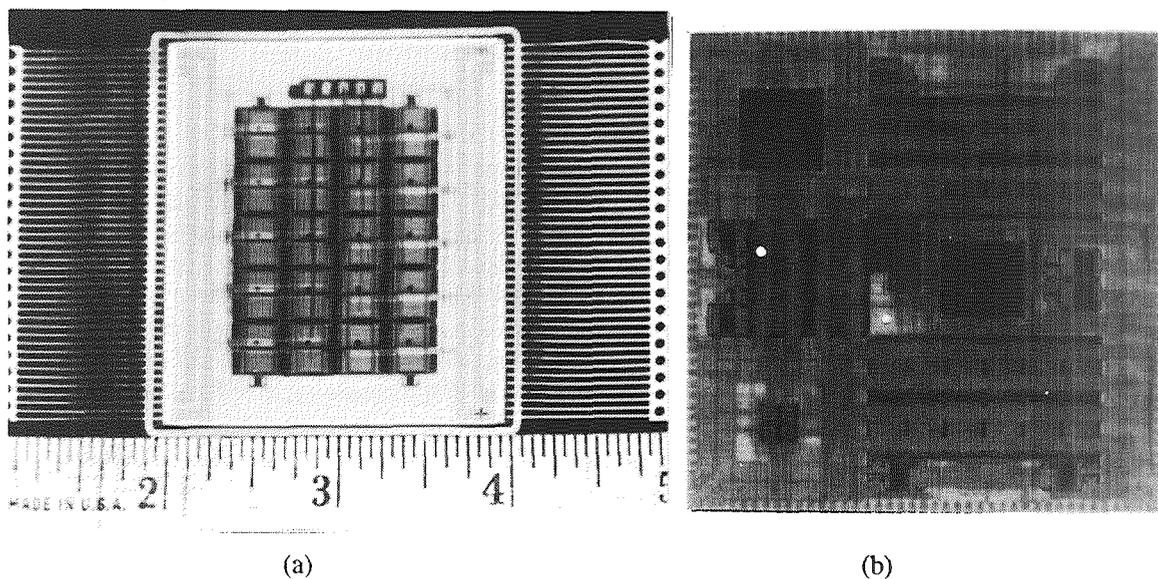


Figure 5. Memory and computer HDI modules. (a) Space-qualified memory module. (b) R3000/R3010 computer core, complete with 256 kilobytes of instruction and data cache.

Future supercomputers will require extraordinary interconnectivity for maximal data bandwidths. Even though monolithic ICs could be constructed to support these high interconnectivity requirements, single chip packages are pin limited, even as they hinder the electrical performance needed by a supercomputer. HDI, particularly 3D-HDI, will allow system designers to re-think and re-define the previously accepted limitations of performance with available packaging technologies. The construction of multi-processor nodes in an exceptionally small space is now possible. Since HDI can support wiring capacities that would allow individual integrated circuits to support thousands of input/output interconnections, multiple-instruction, multi-data stream (MIMD) star topologies for a limited number of 32/64 bit processors operating at very high speeds can be explored. This feature combined with the circuit compression factors provided by 3D-HDI make possible the definition of a new state-of-the-art in computing. Since HDI, for example, allows individual ICs to be placed together within several thousandths of an inch, interconnected with a high performance metal/dielectric system, larger synchronous cache memories, additional co-processors, and associated input/output circuitry can be tightly integrated with a net performance increase, as the new Intel "DX2" microprocessors do at the IC level.

Integrated Sensor/Actuator Assemblies

One trend in advanced packaging is to pack more of a given type of circuitry in the same volume (e.g., more memory). Another is to make a higher performance functional block (e.g., a computing engine). Still another is to integrate more of a system, even a complete system, into a module. Many applications in robotics and automotive electronics could benefit from the ability to build a sensor and processor in the same module. The abilities to place a computer in the arm or joint of a robot or inside the brakes of an automobile allow greatly simplified central control concepts to be applied: less complex communications would be required, bulky wiring harnesses could be reduced or eliminated, and lighter and more rapidly responding systems could be developed. A still more intriguing possibility is that an *entire* electronics system, from sensor to actuator, could be constructed in a two- or three-dimensional HDI system. For example, micro-miniature roving systems, equipped with infrared focal plane arrays or charge-coupled device displays, could be constructed for the exploration of seascapes, the insides of tunnels and mineshafts, or perhaps for the exploration of the surface of another planet. These and other concepts, some still perhaps in the realm of not-so-old science fiction, are achievable through dramatically increased integration capability.

CONCLUSIONS

Through significant research and development activities sponsored by the government, a revolutionary patterned overlay MCM technology has been developed. The HDI technology is as versatile as it is robust, fully capable of meeting the demanding needs of military and commercial applications. It is uniquely capable among MCM technologies in its ability to accommodate a wide variety of existing component technologies, application regimes, operating environments, and ready extension to an even denser three-dimensional form. It is being actively researched for many applications already and is available presently in limited production quantities for evaluation, with large production availability closely following. The benefits of the two- and three-dimensional HDI technology, to reduce the size, weight, and power of electronic systems, and perhaps more importantly, to dramatically improve their performance, make it one of the most significant new technologies for advanced electronics packaging since the original single chip package for the integrated circuit.

REFERENCES

1. Carlson, C.W. *et al.*, "A High Density Copper/Polyimide Overlay Interconnection", *International Electronics Packaging Society (IEPS) Proceedings*, 1988.
2. Bagoklu, H.B. *Circuits, Interconnections, and Packaging for VLSI*, Addison-Wesley, New York, 1990.
3. Cole, H.S. *et al.*, "Polymeric Materials for the GE High-Density Interconnect Process", *Proceedings of 1992 International Society for Hybrid Microcircuits Conference*.
4. Lyke, J.C. *et al.*, "Development and Application of Practical Three-Dimensional Hybrid WSI Technology", *Government Microcircuit Application Conference (GOMAC)*, November 1992.
5. Haller, T.R. *et al.*, "High-Frequency Performance of GE High-Density Interconnect Modules", *Proceedings of the 42nd Electronic Components and Technology Conference, Institute of Electrical and Electronic Engineers*, 1992.

**IMPROVED PERFORMANCE AND SAFETY FOR
HIGH ENERGY BATTERIES THROUGH USE OF
HAZARD ANTICIPATION AND CAPACITY PREDICTION**

Terrill Atwater

U.S. Army Research Laboratory
Electronics and Power Sources Directorate
Power Sources Division

521-33

150491

P-8

ABSTRACT

Prediction of the capacity remaining in used high rate, high energy batteries is important information to the user. Knowledge of capacity remaining in used batteries results in their better utilization. This translates into improved readiness and cost savings due to complete, efficient use. High rate batteries, due to their chemical nature, are highly sensitive to misuse (i.e., Over discharge or very high rate discharge). Battery failure due to misuse or manufacturing defects could be disastrous. Since high rate, high energy batteries are expensive and energetic a reliable method of predicting both failures and remaining energy has been actively sought. Due to concerns over safety the behavior of lithium/sulphur dioxide cells at different temperatures and current drains has been examined. The main thrust of this effort was to determine failure conditions for incorporation in hazard anticipation circuitry. In addition, capacity prediction formulas have been developed from test data. A process that performs continuous, real-time hazard anticipation and capacity prediction has been developed. ¹ The introduction of this process into microchip technology will enable the production of reliable, safe and efficient high energy batteries.

INTRODUCTION

Each year millions of dollars are spent on lithium batteries for use in portable electronics equipment. Because of their superior rate capability and service life over a wide variety of conditions, lithium batteries are the power source of choice for many equipment applications. There is no convenient method of determining the available capacity remaining in partially used lithium batteries; hence, users do not take full advantage of all the available battery energy. In order to maintain readiness, users currently replace batteries on a conservative schedule. This practice results in the waste of millions of dollars in battery energy every year. In addition to the inability to determine the remaining capacity, high rate batteries are highly sensitive to misuse. To preclude this, safety devices are currently included in battery packs. Generally batteries contain three safety components: electrical fuses, thermal fuses and diodes. The electrical fuse protects the cells from sourcing too much current, the thermal fuse protects the battery pack from excessive heat and the diode shields primary cells from being charged. Charging primary cells can have catastrophic results. These devices are passive and provide maximum

protection for all discharge scenarios, limiting the capabilities of the battery by imposing worst case protection for all discharge conditions.

It is a well documented and accepted that the available capacity in a lithium battery is a function of the conditions that the battery has been subjected. Capacity remaining is a complex function of current drain, temperature and time. Therefore a reliable method of predicting remaining capacity has been actively sought. ²⁻⁷ External devices are available for most battery systems. However these devices are, in many cases, not portable and imprecise. Therefore a continuous internal means of determining remaining capacity is desirable. Lithium/sulphur dioxide cell behavior at different temperatures and current drains have been examined. This examination has resulted in the establishment of discharge efficiency formulas. ² Utilization of these formulas has given rise to a capacity prediction algorithm. In addition the safety aspects of electrochemical systems are a priority. ⁸⁻¹⁰ The concern over safety has lead to the incorporation of safety devices into batteries. ⁸ To achieve the full power potential of high energy batteries the development of active safety devices has been examined. Batteries are less tolerant to misuse the deeper they are discharged. Knowledge obtained through safety testing incorporated into microelectronics technology allowed the development of active safety devices. These devices designed to anticipate and detect hazardous conditions allow for increased power loads to be placed on the battery without the fear of failure.

EXPERIMENTAL

Lithium/sulphur dioxide cells, produced under government contract, were discharged at various temperatures and discharge rates. Cell discharge included typical as well as abusive conditions. Tests were temperature controlled using a Blue M, Model 1004-3B environmental chamber. Discharge rates were controlled using an internally fabricated constant current electronic load. This electronic load was designed to be a constant current load without forced discharge or external power supply aid. Test conditions were continually monitored and controlled with an ACROSYSTEMS Acro-400 data acquisition and control unit in conjunction with a personal computer. Acquisition and control software was internally generated. Temperature, discharge rate and cell voltage were continually recorded for analyses.

Typical discharge rates ranged from 0.02 mAmp/cm² to 10 mAmp/cm². This normalized discharge rate allowed for different cell sizes to be compared. The tests were conducted in temperatures between -40°C and 70°C. Cells were subjected to constant current discharge, pulse discharge and intermittent discharge. During pulse and intermittent discharge, conditions were varied as well as kept constant. Cell discharge capacity was determined at a two volt per cell cut off. The lithium/sulphur dioxide electrochemical system has an open circuit voltage of three volts. Cell discharge efficiency, determined by the ratio of delivered capacity to theoretical capacity, is graphically represented in Figure 1. ² This set of curves can be described by a surface represented by Equation 1.

Discharge Efficiency of lithium/sulphur dioxide cells

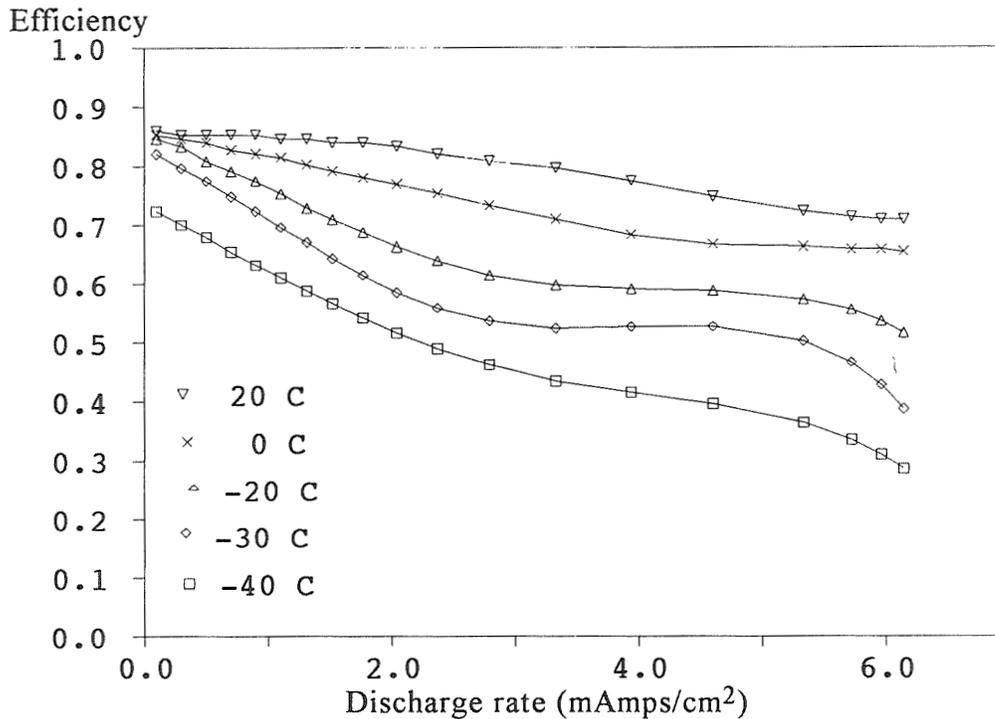


Figure 1. Graphic representation of relationship between discharge efficiency, temperature and constant current discharge rate for lithium/sulphur dioxide cells.

$$Ed = \sum_{x=0}^3 \sum_{y=0}^3 C_{x,y} \cdot i^x \cdot T^y$$

Equation 1. Polynomial fit of the surface relating discharge efficiency (Ed) to discharge current (i) and temperature (T).

Charging, forced over discharge and high current discharge are examples of abusive discharge conditions. These conditions generally lead to cell failure and were performed during the safety analysis.¹⁰ Typical cell failure involves the opening of the cell's vent. In addition partially discharged cells were interrogated using DC pulse measurement and AC impedance spectroscopy using an impedance spectroscopy measurement system. The impedance spectroscopy measurement system included a Hewlett Packard HP9836 computer, a Schlumberger 1286 Electrochemical Interface and a Schlumberger 1250 Frequency Response Analyzer.³ The effects of discharge rate, depth of discharge and temperature on cells were observed through the response to this stimulation. Data obtained during all tests were used to develop a capacity prediction algorithm as well as hazard anticipation circuitry.

DESIGN

The remaining capacity of batteries is often determined by a simple technique that monitors the coulombic drain (Amperes-seconds) from the battery. This information tells the user the capacity removed from the battery. When this information is subtracted from nominal capacity, remaining capacity can be determined. This technique is valid when discharge conditions remain constant. Variation in battery efficiency due to discharge conditions can result in large errors. Since discharge conditions of a battery are continually changing, it is desirable to utilize battery discharge efficiency on a continual basis. Designed capacity prediction circuitry calculates effective capacity removed from a battery based on discharge efficiency on a continual basis. In order to calculate effective capacity removed discharge rate, temperature and time must be determined. ²

The voltage drop across a resistor is used as a current sensor. Unfortunately, there is valuable power wasted in the resistor when this approach is used. Another approach to monitor the battery's current drain is to use one of the components already incorporated in the battery. Using the diode as the sensor is not recommended because connecting anything across the diode would prevent the diode from performing its protective function. The voltage drop across a fuse is a function of the current passing through it and its ambient temperature. However variation in the fuse can cause calculation errors. Therefore a resistor with a value as low as possible is used.

The voltage across the resistor will be relatively small and must be amplified in order to bring it up to levels that can be processed. In addition the voltage drop will fluctuate, due to varying discharge loads. Therefore to maintain accuracy it is extremely important to continuously monitor the voltage across the resistor and store this information. One method of capturing this voltage is to use an integrator that would sum up the voltage over time. This summation represents an average current therefore the integration time should be minimal. The integrator can be of either analog or digital design; however, an analog integrator is far simpler to implement. An inexpensive method of integrating a voltage is to use a pair of PNP transistors and configure them as a conventional differential amplifier. A constant current source is needed to prevent changes in cell voltage from effecting the data being collected. Besides amplifying the resistor's voltage, the differential amplifier converts the voltage to a current. This current is used to charge a capacitor. The voltage developed across the capacitor represents capacity removed and is sensed by a comparator. The comparator and processor converts the analog voltage to digital data. ⁴

The capacitor and resistor used in the circuit are tuned such that a threshold voltage represents one coulomb capacity removed. The capacitor is discharged by a switch controlled by the processor whenever the threshold voltage is achieved. This cycle is repeated each time one coulomb capacity is removed from the battery. An internal clock is used to record time for each cycle. Temperature data is acquired similarly, however

amplification is not necessary. A current source whose output is a function of temperature is commercially available and inexpensive. ⁴

The processor calculates current by calculating the ratio of one coulomb and the time required to perform the integration cycle. With current and temperature information the efficiency of the reaction can be calculated. Adjusted capacity removed (one coulomb/efficiency of reaction) is subtracted from the previous capacity remaining, this value represents the current capacity available. For processors with limited computing power a look up table can be utilized without adding significant error. ^{1,4}

A control feature added to the processor enables hazard anticipation to be performed. This control requires the addition of voltage sensors across each cell and control elements within the battery. Example control elements include: micro controlled current suppressers and electronic fuses. A routine utilizing capacity remaining, temperature, discharge current and cell voltage is used to control the use of the battery. Examples of typical control include:

- 1) Allowable discharge current reduced as available discharge capacity decreases.
- 2) Load disconnected from battery when internal temperature reaches a preset level, allowing the battery to cool off.
- 3) Load permanently disconnected from load when the voltage of one cell in the battery drops below a preset level.

A simple version of battery control uses passive safety devices and an active safety device that permanently disconnects the battery from any load when hazardous conditions are detected.

Incorporated in the design is the ability to output to the user the information gathered and stored by the processor. A light emitting diode readout positioned on the battery will display to the user the remaining capacity in the battery. In order to conserve battery energy this readout is activated by a switch located on the battery surface. In addition a continual digital readout using the battery connector will output data to equipment. This output will allow for a real-time indication of the abilities of the batteries without removing the battery from the equipment. This output requires the equipment to have the ability to read and decipher the signals being outputted from the battery. ¹

Battery data output is achieved using a NMOS FET open drain with the source grounded. The serial data output has one start bit followed by eight data bits. The most significant bit of the data will be first and the data will represent a two character hexadecimal word. This hexadecimal word is used as a code that represents battery condition. Table 1 contains a listing of data codes and their meanings. In the output logical "0" is represented by high voltage and logical "1" is represented by low voltage. The data output will be provided every second at a bit transfer rate, in the order of, 800 bits per second. ¹

HEX CODE	INDICATION
00	00 percent capacity remaining in battery
01	01 percent capacity remaining in battery
...	...
0A	10 percent capacity remaining in battery
...	...
64	100 percent capacity remaining in battery
...	...
70	User activated shutdown
71	Battery group 1 cell 1 failure
72	Battery group 1 cell 2 failure
...	...
7A	Battery group 2 cell 5 failure
7B	Battery group 1 fuse activation
7C	Battery group 2 fuse activation
...	...

Table 1. Partial listing of output codes for a ten cell battery pack configured in two five cell strings.

Processor controlled circuitry that incorporates hazard anticipation, remaining capacity indication and data output has been designed. This circuit contains hardware and software for determining remaining capacity and control. The hardware is controlled by software, this makes the design suitable for different electrochemical systems and battery configurations.

RESULTS AND DISCUSSION

Breadboard circuits were fabricated to test the validity, accuracy and operation of the design. The breadboard allowed for quick software changes to accommodate different cell configurations. The breadboards incorporated both hazard anticipation and passive safety devices. Operation tests were performed using 'D' size and '1/3 C' size lithium/sulphur dioxide cells. Cell configuration included: Two five cell strings discharged in parallel and one, five and ten cell strings discharged in series. During discharge temperature, voltage and current were monitored independent of the circuitry. Discharge conditions for the validity checks were the same as described in the experimental section of this paper. Discharges were halted at random depth of discharge.

To evaluate the operation of the hazard anticipation, single cell and multi-cell batteries were stressed with known failure modes. Tests were considered successful when the circuit disconnected the battery and no cells within the battery vented. Evaluation of the capacity prediction algorithm was performed by recording the predicted remaining capacity of partially discharged batteries. The batteries were then discharge at a rate of 5

mAmps/cm² at 25°C. The capacity delivered was then compared to the capacity of a fresh battery at this discharge condition. The fresh battery capacity was determined by a twenty battery average. The results of this comparison is shown in Figure 2. Figure 2 shows a prediction error of ±5% with a median value slightly skewed to the negative. A negative error denotes a low capacity remaining prediction. ²

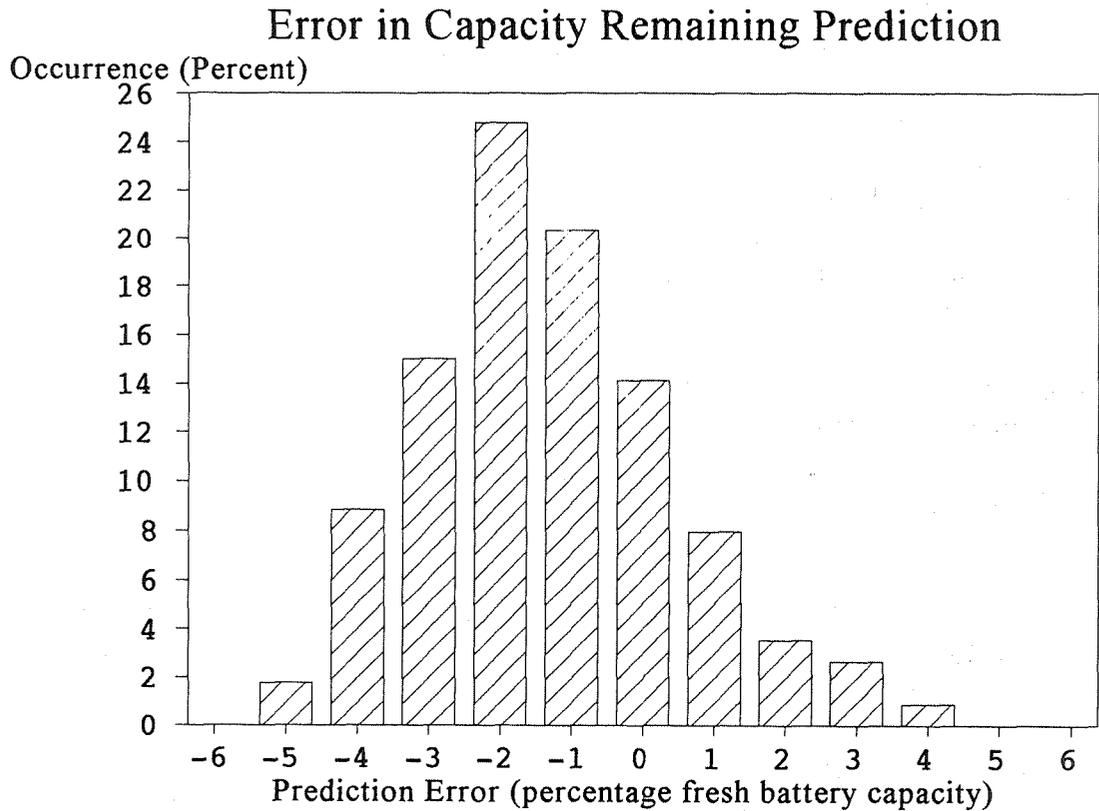


Figure 2. Distribution of error in capacity remaining prediction.

Data output capability was performed using a developmental radio. Tests showed that the data supplied from the battery was able to be read and interpreted by the radio. The radio was then capable of displaying this information to the user. These results combined with the low capacity remaining prediction error and successful hazard anticipation results shows that the incorporation of hazard anticipation and capacity prediction provides the means for more efficient utilization of battery energy by the user.

REFERENCES

1. T. Atwater, A. Bard, B. Testa and W. Shader, 'Smart Battery Controller for Lithium Sulfur Dioxide Batteries', U.S. Army, LABCOM, Research and Development Technical Report, SLCET-TR-91-36, August 1992
2. Terrill Atwater, 'Performance Prediction for Military Lithium Sulphur Dioxide Batteries', Proceedings 180th Electrochemical Society Meeting, October 1991.
3. T. Atwater, C. Kanaris and K. Gaetano, 'Improved Impedance Spectroscopy Technique for Status Determination of Production Li/SO₂ Batteries', Proceedings 35th Power Sources Symposium, 1992.
4. T. Atwater, A. Bard, 'An Universal Inexpensive Battery Effective Capacity Indicator', in course of publication.
5. M. L. Gopikanth, 'State-of-Charge Indicators for Aqueous Primary Batteries — A Review', Proceedings 34th Power Sources Symposium, 1990.
6. F. Walsh, 'Determination of State-of-Charge in Li/SOCl₂ Cells', Proceedings 34th Power Sources Symposium, 1990.
7. G. DiMasi, T. Atwater, M. Brundage and L. Jarvis, 'Temperature Performance Characteristics of Lithium-Thionyl Chloride Cells with Various Cathode Additives', Proceedings 170th Electrochemical Society Meeting, October 1986.
8. M. Brundage, G. DiMasi, L. Jarvis and T. Atwater, 'Significant Advances in the Safety and Technology of Lithium-Sulfur Dioxide Batteries', Proceedings 32nd Power Sources Symposium, 1986.
9. G. DiMasi and J. Christopoulos, 'The Effects of the Electrochemical Design Upon Safety and Performance of The Lithium-Sulfur Dioxide Cells', Proceedings 28th Power Sources Symposium, 1978.
10. G. DiMasi, 'Behavior of Li/SO₂ Cells Under Forced Discharge', Proceedings 27th Power Sources Symposium, 1976.
11. David Linden Ed., 'Handbook of Batteries and Fuel Cells', McGraw-Hill Book Company, 1984

THIN RECHARGEABLE BATTERIES FOR CMOS SRAM MEMORY PROTECTION

Dennis N. Crouse
EIC Laboratories, Inc.
111 Downey Street
Norwood, Massachusetts 02062

522-33
N 93-25583
150492

P-6

ABSTRACT

New rechargeable battery technology is described and compared with classical primary battery back-up of SRAM PC cards. Thin solid polymer electrolyte cells with the thickness of TSOP memory components (1 mm nominal, 1.1 mm max.) and capacities of 14 mAh/cm² can replace coin cells. SRAM PC cards with permanently installed rechargeable cells and optional electrochromic low battery voltage indicators will free the periodic PC card user from having to "feed" their PC cards with coin cells and will allow a quick visual check of stored cards for their battery voltage status.

INTRODUCTION

The advantages of CMOS SRAM's fast access time, access time equal to cycle time, equal duration read/write cycles and low voltage write capability are tempered by SRAM's requirement of a small data retention current when disconnected from the host computer's power. For this reason PC cards with SRAM require a battery back-up in order to retain their memory when disconnected from system power. The PCMCIA standard specifies a battery location at the end of the PC card opposite the connector (Figure 1) and PC card manufacturers use a coin type battery in a removable battery holder (Figure 2) in this location. The PCMCIA standard section 4.12 also makes provision for the PC card to compare the battery voltage with two reference voltages and provide status signals (BVD1 · BVD2 = 0) when the battery needs to be replaced.

PRIMARY COIN CELLS

The low profile of the 2025 coin cell makes it the cell of choice for memory cards. The 20 represents its diameter in millimeters and the 25 represents its thickness in tenths of a millimeter. The voltage and capacity of this cell varies depending upon the chemistry used in the cell. Two types of three volt coin cells are commercially available: the CR2025 uses lithium - manganese dioxide chemistry and has a capacity of 140 mAh, while the BR2025 uses lithium - polycarbon monofluoride chemistry and has 120 mAh of capacity.

Battery capacity can be translated into PC card memory life if both the SRAM's data retention supply current (ICCDR) and average storage temperature are specified. Figure 3 is a plot of the battery capacity as a function of average storage temperature for 1 Mbyte of SRAM (eight NEC uPD431000A). The two sloping lines indicate the capacity for one month and two years of data retention. As indicated on this plot a CR2025 battery will provide 2 years/MByte of data retention at 25°C and only 1 month/MByte at 80°C. This steep temperature dependence puts the responsibility for memory retention upon the user, as they must keep their cards cool and periodically insert them in a computer to check on battery status.

RECHARGEABLE TLOP CELLS

Recent advances in thin solid polymer electrolytes have made it possible to provide rechargeable cells in the same thickness as SRAM in thin small-outline packages (TSOP). The TSOP SRAM have a height above the circuit board of only 1.0 mm nominal and 1.1 mm maximum. Since the rechargeable cells are in thin large-outline packages (2 to 30 cm²), we call them TLOP cells, an example of which is shown in Figure 4. The capacity of these TLOP rechargeables depends upon both their surface area and chemistry. Figure 5 is a plot of battery capacity per square centimeter versus cell thickness for several commercially available coin and rectangular primary cells as compared with a solid polymer electrolyte rechargeable TLOP cell.

Comparing the two 1 mm thick cells in Figure 5 shows that the rechargeable LIMO TLOP cell has capacity per square centimeter equal to commercially available primary coin and rectangular cells of comparable thickness.

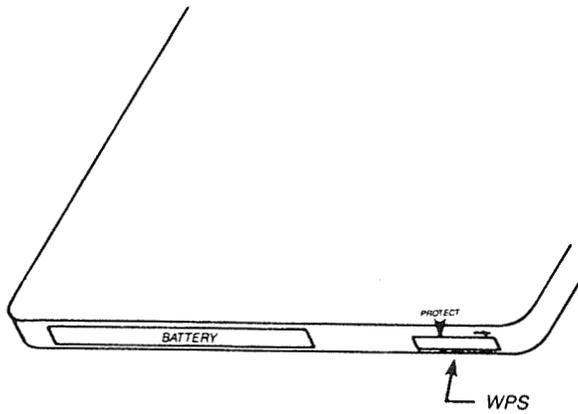


Figure 1. Battery and Write Protect Switch (WPS) Location on PC Card.

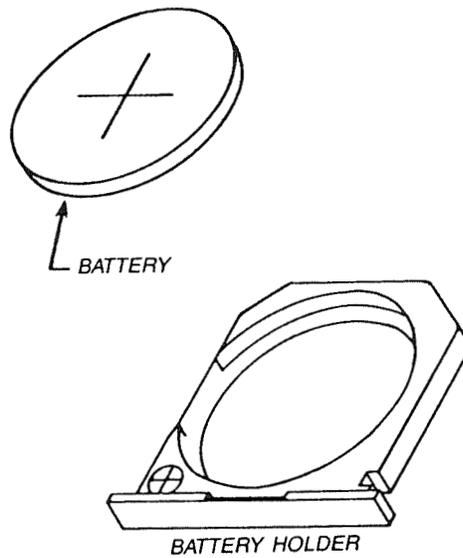


Figure 2. Coin Cell and Battery Holder.

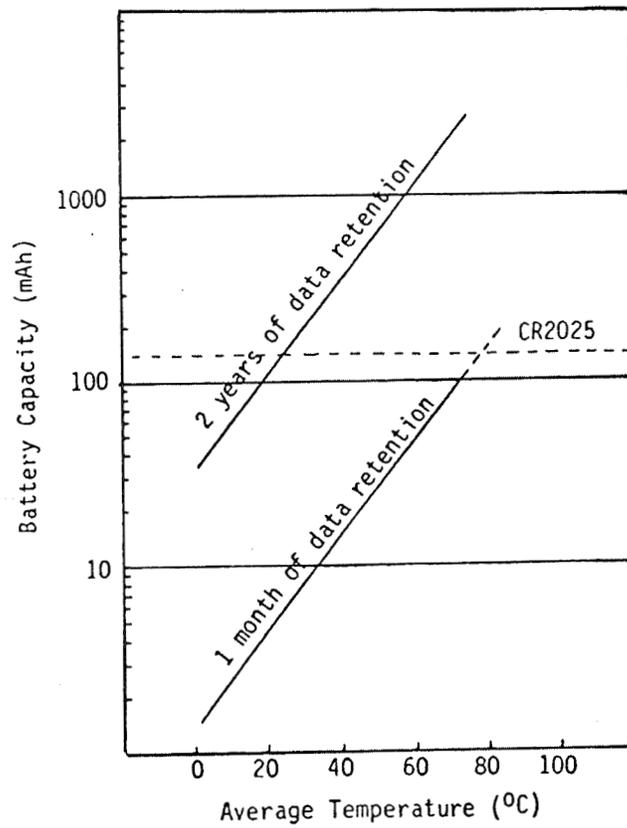


Figure 3. Battery Capacity Required for Data Retention in Eight 1 MBit SRAMs (based upon NEC uPD 431000A).

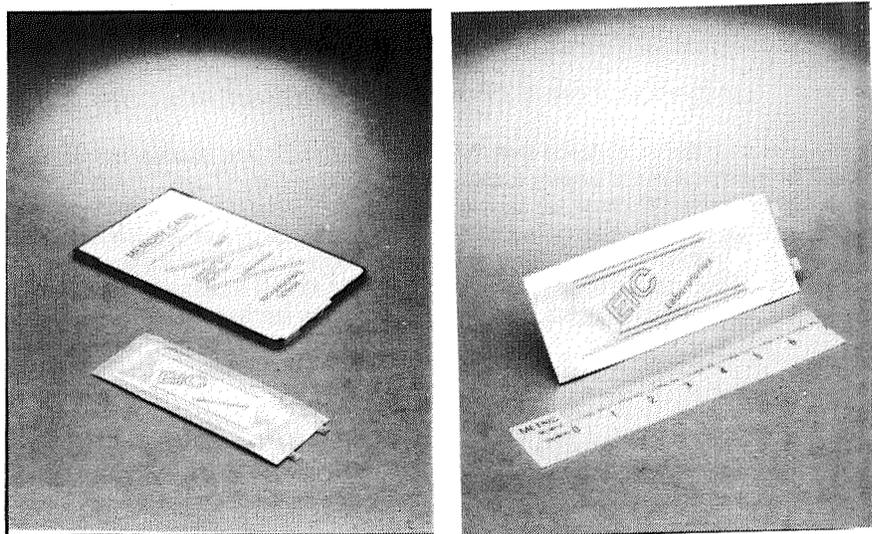


Figure 4. Rechargeable TLOP Cells for PC cards.

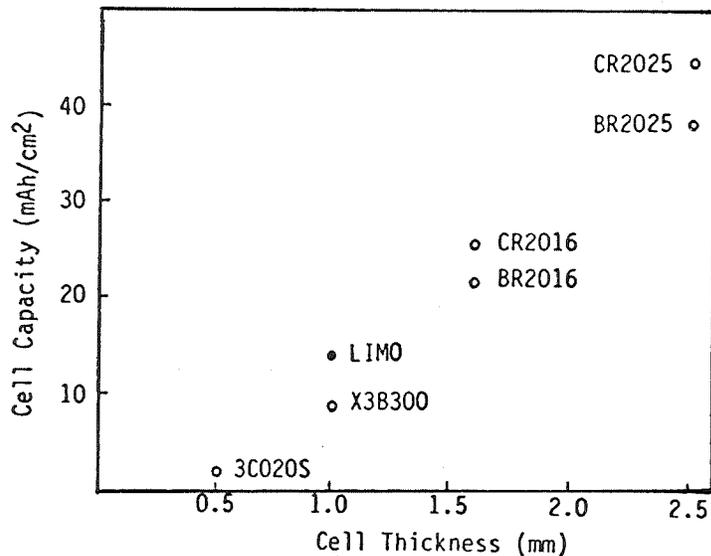


Figure 5. Nominal Cell Capacity per Square Centimeter Versus Cell Thickness for Primary (○) and Rechargeable (●) Cell.

The LIMO TLOP cell has a theoretical capacity of 14 mAh/cm² and can deliver 90% of this capacity at greater than 2 volts over a wide temperature range with loads equal to thirty two TSOP SRAMs. Recharge of the LIMO cell takes place from 3.2 to 3.7 volts and these cells have been charged and discharged more than 100 times with retention of greater than 50% of their initial capacity. Figure 6 is an overlapping series of charge/discharge plots of a LIMO cell. Figure 7 is a cross-sectional view of a LIMO cell.

RECHARGEABLE TLOP CELLS IN PC CARDS

SRAM PC cards with rechargeable TLOP cells, such as the LIMO cell, are capable of being recharged, while in use, from the host computer's power. This results in SRAM PC cards in which battery back-up is automatically maintained from the perspective of the periodic user. The TLOP cell's length and width can be customized for specific PC card applications. For instance a 4.5 by 7 cm TLOP cell could be surface mounted on one side of a PC card's circuit board and sixteen TSOP SRAM could occupy the other side. If this TLOP cell were a LIMO cell it would have a capacity in excess of 420 mAh or more than three times a CR2025 coin cell.

OPTIONAL LOW BATTERY VOLTAGE INDICATOR

To provide the user with the confidence that the battery has sufficient charge, an electrochromic low battery voltage indicator could be mounted in the battery location on the end of the PC card. Prototypes of voltage indicators installed in commercially available PC cards are shown in Figure 8. The indicator would also allow stored PC cards to be visually scanned to insure data retention in a fraction of the time required to sequentially insert the PC cards into a host computer for battery status verification. Electrochromic displays offer the advantage of a nonvolatile display requiring no refresh current. This combination of a rechargeable TLOP cell and an electrochromic low battery voltage indicator will result in SRAM PC cards being entrusted with information storage in applications involving only periodic updating or editing.

CONCLUSION

Solid electrolyte technology now exists which makes possible the manufacture of rechargeable cells with the thickness of TSOP components. These thin large-outline cells, called TLOP cells, can be permanently installed in SRAM PC cards and will free the periodic PC card user from having to "feed" their PC cards with coin cells.

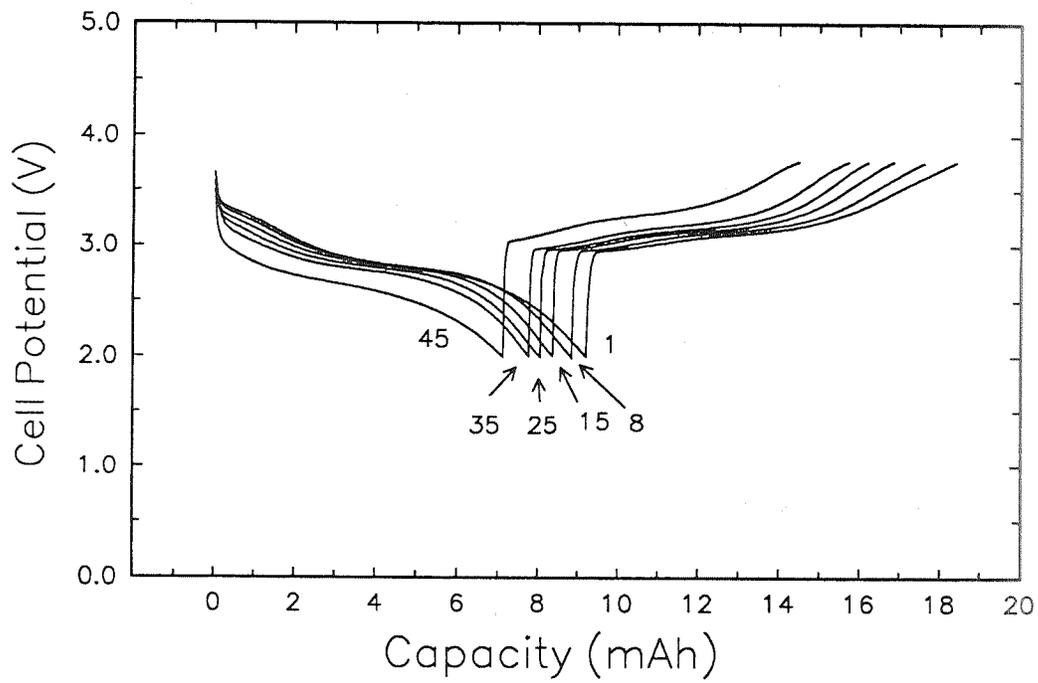


Figure 6. Overlapping Charge/discharge Plots of a LIMO Cell with Cycle Number Indicated.

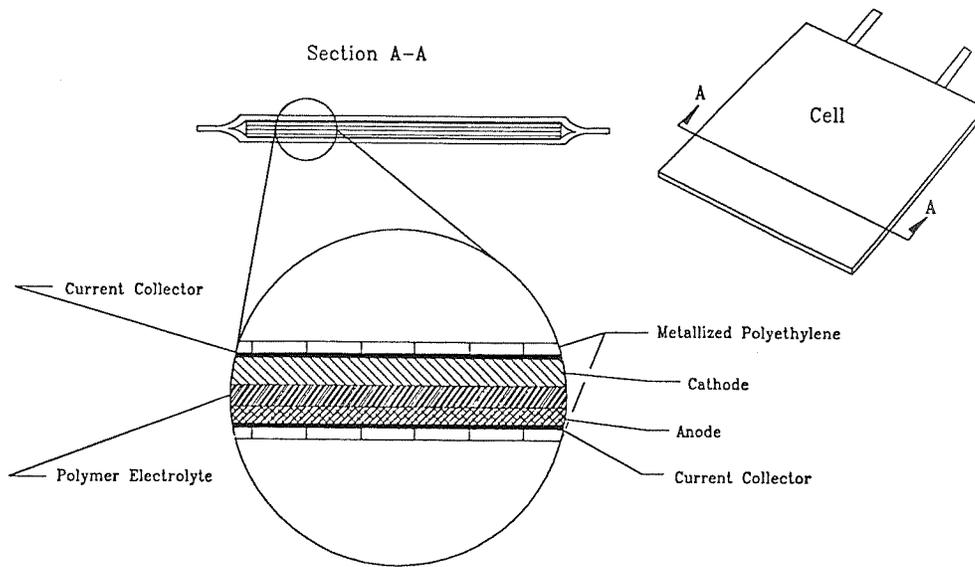


Figure 7. Cross-sectional View of LIMO Cell.

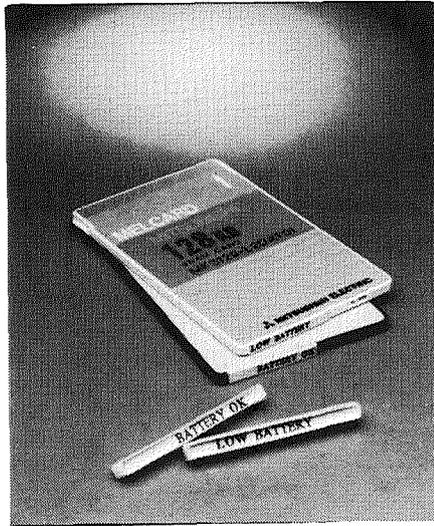


Figure 8. Prototype voltage indicators assembled and installed by EIC Laboratories.

amt

PASSIVE STACKING FOR IMPROVED VIBRATION ISOLATION

This paper was not submitted for inclusion in these proceedings. You may wish to purchase the audiocassette of this presentation by contacting:

**The Technology Utilization Foundation
41 East 42nd Street, Suite 921
New York, NY 10017
Ph.: (212) 490-3999**

If you would like further information on this presentation, please contact:

**David A. Noever
AST, Biophysics Branch
NASA Marshall Space Flight Center
Space Sciences Laboratory
Mail Code ES76
MSFC, AL 35812
Ph.: (205) 544-7783**

omit

**ADVANCED MATERIALS PART 2:
CERAMICS AND COMPOSITES**

PRECEDING PAGE BLANK NOT FILMED

221.

290
~~INTENTIONALLY BLANK~~



323-76
N 93-25584
594584
p-8

**A Novel Method for Characterization of Superconductors:
Physical Measurements and Modeling of Thin Films**

**B. F. Kim, K. Moorjani, T. E. Phillips
F. J. Adrian, J. Bohandy and Q. E. Dolecek**

**The Johns Hopkins University
Applied Physics Laboratory
Johns Hopkins Road
Laurel, Maryland 20723**

ABSTRACT

A method for characterization of granular superconducting thin films has been developed which encompasses both the morphological state of the sample and its fabrication process parameters. The broad scope of this technique is due to synergism between experimental measurements and their interpretation using numerical simulations. Two novel technologies form the substance of this system: the magnetically modulated resistance method for characterizing superconductors, and a powerful new computer peripheral, the Parallel Information Processor card, which provides enhanced computing capability for PC computers. This enhancement allows PC computers to operate at speeds approaching that of supercomputers making atomic scale simulations possible on low cost machines. The present development of this system involves the integration of these two technologies using meso-scale simulations of thin film growth. A future stage of development will incorporate atomic scale modeling.

INTRODUCTION

In this paper, we describe a system for characterization of superconductor thin films which encompasses both physical properties of thin films and parameters which control their fabrication. This tool is naturally suited for development of thin film superconductor materials and for quality control. A comprehensive assessment is achieved by physical measurements and computational simulation of the measurement responses. This system is implemented by the integration of two novel technologies: the magnetically modulated resistance (MMR) method for characterization of superconductors,¹⁻⁵ and the parallel information processor (PIP)⁶⁻⁸ for desktop supercomputing with PC computers.

The MMR method measures the magnetic field derivative of resistance as a function of temperature. This measurement has been shown to be sensitive to effects which occur in granular superconductors.² Granular superconductors consist of ensembles of small single crystals (grains) which are in contact with one another. Two grains separated by a common boundary constitute a superconductor structure called a weak link which has superconductor properties different in some respects from that of a single grain.⁹ Granular superconductors, therefore, exhibit effects that are attributable to the weak links which constitute their morphology. It is important to account for these effects since virtually all high temperature superconductors, for practical purposes, are granular.

Our primary interest here is the characterization of granular thin film superconductors which are of paramount importance in superconductor devices. In the following discussion we will describe briefly how the MMR measurement characterizes superconductors. We then discuss briefly how the MMR measurement can be related to thin film fabrication parameters by computational simulation of the MMR response signal. Numerical simulation is becoming increasingly important in materials research. In this work, it provides an essential link between physical measurements, and sample morphology and process parameters. Calculations of this type require powerful computing machinery. Atomic scale materials simulations often require supercomputers. The PIP technology mentioned above allows the simulations in this project to be done on a PC. We will discuss this new technology, briefly, and show how the physical MMR measurements and computational elements are being integrated into a system for characterizing superconductor thin films. Finally we will indicate how we intend to extend this system in the future.

THE MMR METHOD

The essential features of the MMR technique are shown in the diagram in Figure 1. The superconductor sample is contained in a variable temperature bath. The sample is subjected to a magnetic field consisting of an ac

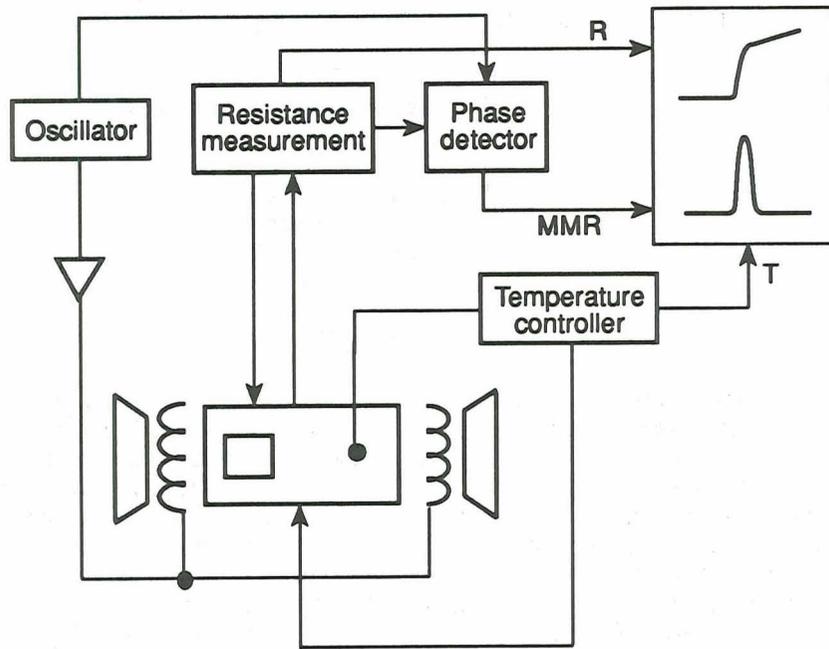


Figure 1. Schematic diagram for an apparatus to measure magnetically modulated resistance.

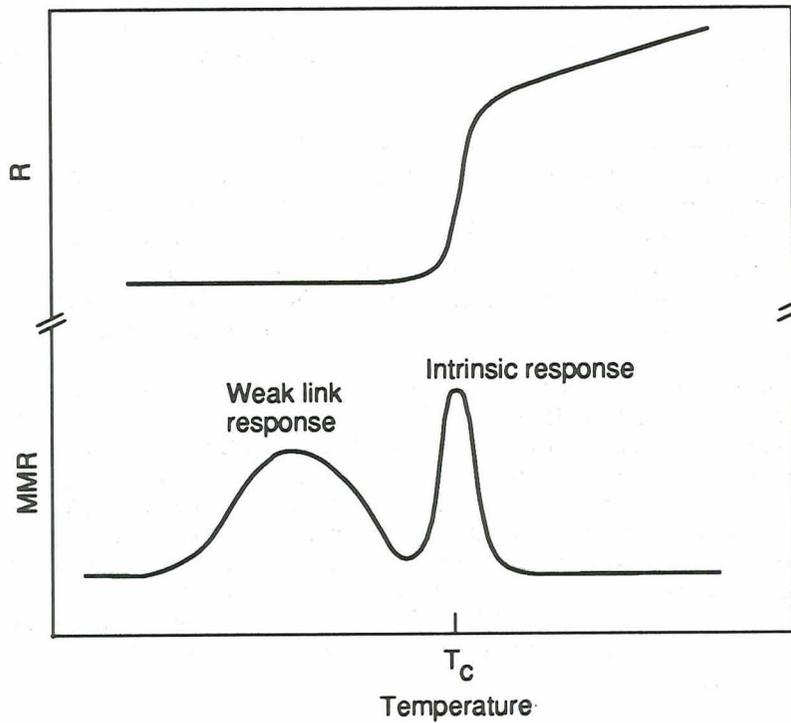


Figure 2. Typical resistance (R) and magnetically modulated resistance (MMR) response curves for a granular superconductor.

component, H_{ac} , and a collinear dc component H_{dc} where $H_{ac} < H_{dc}$. The resistance of the sample is measured at the field modulation frequency by phase detection. This signal, which constitutes the MMR response, is measured as a function of temperature. The method for measuring resistance can be any convenient method. We use four point probe dc resistance measurements and X-band microwave surface resistance measurements in our laboratory. (A detailed description of this method is contained in Reference [10]).

A typical MMR response signal for a granular superconductor is shown schematically in Figure 2. Also shown for comparison is the unmodulated resistance which is measured simultaneously. The MMR response contains two features; a peak at T_c which signifies the superconductor phase transition for individual superconductor grains, and a second peak below T_c which is due to superconductor weak link phase transitions. The relative position and shape of the MMR weak link feature depends upon the current density in the sample, the externally applied magnetic field, and the distribution of grain sizes and weak link coupling strengths in the sample. The dependence of the MMR weak link peak on grain size distribution provides the basis for characterization of the morphology of granular superconductors. In the following section, we briefly describe how the MMR signal is related to the sample grain size distribution.

THE MMR SIGNAL

In this section, we outline the physical interpretation of the MMR signal without detailed discussion of the underlying physics. The reader who is interested in more detail is referred to Reference [11]. An expression for the MMR signal vs. temperature was inferred under the assumptions that the weak links are Josephson junctions⁹ with equal coupling strength per unit area, and that the applied current is unidirectional in the sample. Under these conditions it can be shown that the MMR signal has the form

$$MMR(T) = KWR^2 \frac{e^A}{I} F(T) \frac{d}{dB} \sum_L \frac{\sin(\pi B/B_o)}{\pi B/B_o} f(L) \quad (1)$$

where

$$R = \frac{1}{1 + e^A}$$

$$A = -\left(1 - \frac{I_c}{I}\right)$$

$$I_c = \sum_L \frac{WK}{B_o} F(T) \frac{\sin \pi B/B_o}{\pi B/B_o} f(L)$$

$$B_o = \frac{K}{L}$$

In the expression above, K is a constant, I_c and I are the critical and applied currents, $F(T)$ is the temperature dependence of I_c , B is the applied magnetic field, B_o the critical field for a junction of length L , W and L are the width and length of a junction, and $f(L)$ is the distribution of junction lengths. The function $F(T)$ is known but is not related to the sample morphology. The junction length is the length along a grain boundary which is perpendicular to the current, and which is bounded by grain boundaries parallel to the current as illustrated by Figure 3. (For our purposes, all grain boundaries are resolved into directions parallel and perpendicular to the applied current.) The junction size distribution $f(L)$ is related to the grain size distribution $g(L_g)$ in a complex way, but this relationship can be calculated by a method of statistical inference. An example of this relationship is illustrated in Figure 4 for a Gaussian grain size distribution. Thus from Equation (1) and the relation between $f(L)$ and $g(L_g)$, we are able to relate the MMR signal to the grain size distribution of the sample.

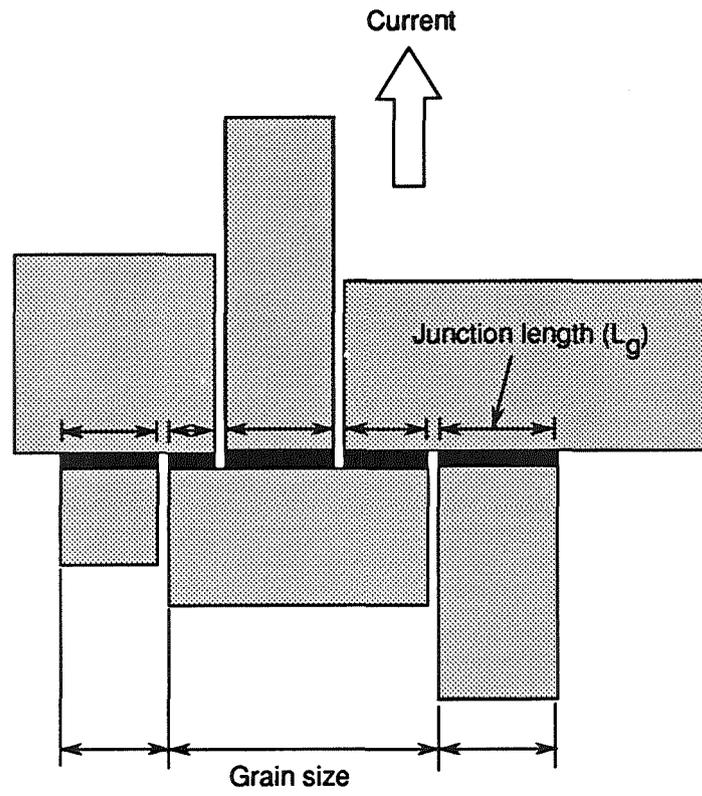


Figure 3. An illustration which shows the relation between grain size (L) and junction size (L_g) in an ensemble of closely packed grains.

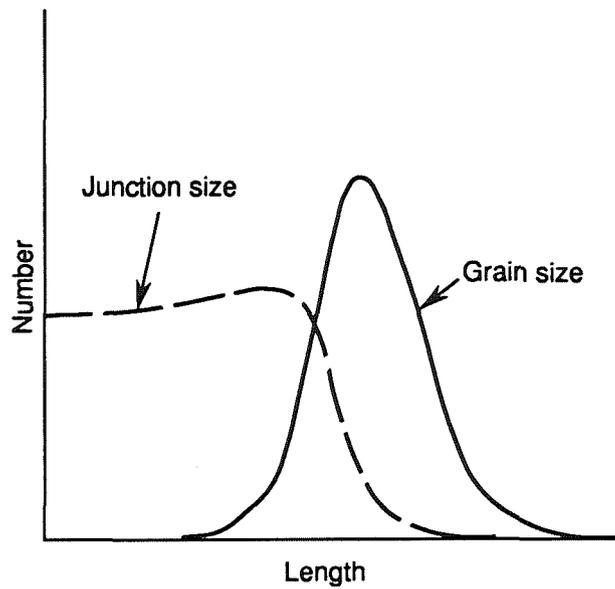


Figure 4. An example of the junction size distribution which results from a Gaussian grain size distribution.

The simulation of the MMR signal described here provides a way, in principle, to assess thin film morphology which is, for our purposes, the grain size distribution of the thin film. In actual practice, one cannot directly solve Equation 1 to obtain $f(L)$. Consequently, one must essentially estimate the function $g(L_g)$ and then compare the calculated MMR response with the experimental response to gauge the accuracy of the estimate. The problem, then, is how to estimate $g(L_g)$. Because the distribution $g(L_g)$ is directly related to the conditions which control the growth of the thin film, our approach is to estimate $g(L_g)$ by modeling the growth of thin films. This provides a way to obtain the grain size distribution from MMR measurements and at the same time extends the application of the MMR method to control of process parameters of thin film growth. In the next section, we discuss our work in modeling thin film growth to obtain estimates of $g(L_g)$.

SIMULATION OF THIN FILM GROWTH

The simulations of vapor phase thin film growth described here, generally apply to state of the art epitaxial thin films such as, for example, $\text{YBa}_2\text{Cu}_3\text{O}_{7-y}$ on LaAlO_3 substrates. These types of films consist of grains of various sizes which have a common orientation normal to the substrate. Growth occurs as a result of the vapor phase deposition of the constituent species, which constitute the superconductor, onto the substrate. The vapor phase constituents are commonly created by rf sputtering or laser ablation from a bulk superconductor target. The vapor initially condenses on the substrate and forms a number of seed crystals of the superconductor species. The seed crystals grow until their boundaries encounter neighboring crystals thereby forming common grain boundaries between the neighboring crystals. The size of a particular crystal is then determined by the presence of neighboring crystals which limits their lateral growth to the area encompassed by their grain boundaries.

The grain size distribution is inferred by numerical simulation of this mechanism of thin film growth. The simulation is presently implemented as a meso-scale model with two adjustable parameters. These are the probability P per unit time that a seed crystal will form at any given location on the substrate, and the rate of growth, R , of a crystal on the substrate. The simulation proceeds by determining the location of new seed crystals on the substrate in each unit of time according to the parameter P , and growing each existing crystal on the substrate according to the parameter R . This process continues until there are no voids on the substrate. The resulting grain size distribution is obtained directly by assessing the area of each crystal (grain). Thus, this simulation provides the shape of $g(L_g)$ which is determined by the two parameters P and R .

We plan to extend the film growth simulation using atomic scale simulation of the deposition. The formation of seed crystals and crystal growth will be simulated using interatomic interactions and the physical parameters of the growth process such as substrate temperature, oxygen pressure, and certain parameters of the reactant species. This will permit calculation of the meso-scale parameters P and R . The atomic scale model will be appended to the meso-scale model to provide the grain size distribution which will then be related to physical process variables.

AN MMR CHARACTERIZATION SYSTEM

The concepts discussed in the preceding sections are being integrated into a system for characterization of superconductor thin films. The initial version of this system, shown diagrammatically in Figure 5, is based upon a comparison between experimental MMR measurements and numerical simulations of MMR response. The experimental measurements are obtained separately and stored in a computer file for subsequent off-line comparison with simulation. The simulation begins with the meso-scale thin film growth model which calculates the grain size distribution. The grain seed probability and growth rate parameters, P and R , are initially estimates which are entered manually. The junction size distribution, calculated from the grain size distribution, is used to calculate the MMR response. (The MMR response also depends upon some parameters not shown in Figure 5; e.g., magnetic field and current.) Comparison of the simulated and experimental MMR responses yields an MMR error from which new values of P and R are inferred. The process iterates until the MMR error falls within predetermined bounds. In this way the grain size distribution of the thin film is ascertained.

The parameters P and R are also determined in this system, although this knowledge is primarily of academic interest. In order to relate this to more fundamental physical parameters we intend to extend this system using atomic scale modeling of the deposition process. Figure 6 shows this extended system. The P and R parameters are calculated in the atomic scale model, as discussed previously, using process variables such as substrate temperature, oxygen pressure, and certain parameters of the reactants (e.g., momentum, excitation, etc.). Thus, in this system, the MMR response will be related to the thin film process parameters.

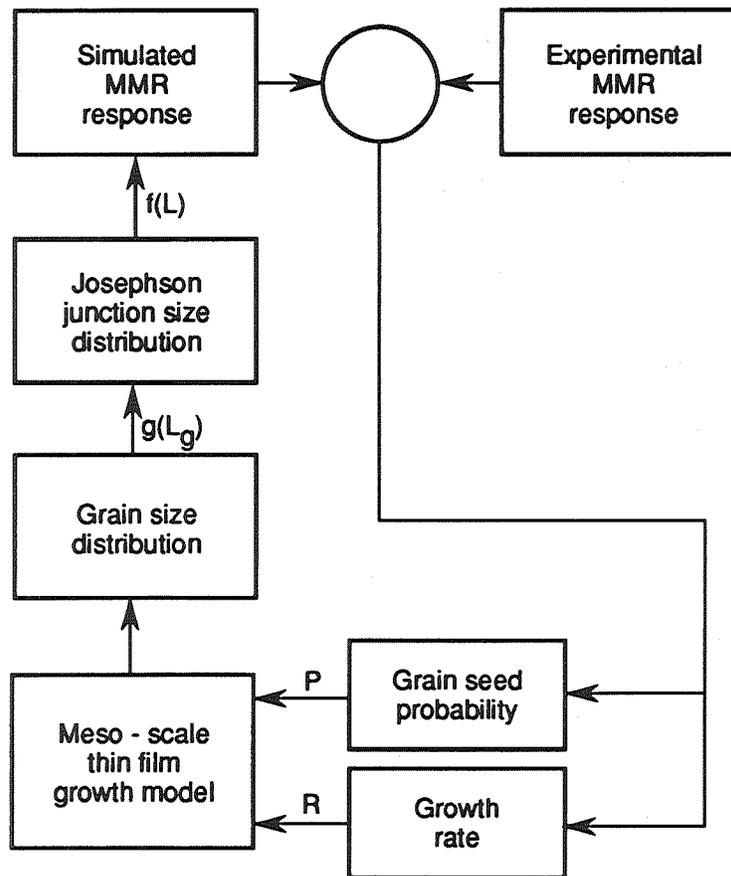


Figure 5. Schematic diagram of an MMR characterization system which uses meso-scale modeling.

The systems described above, particularly the extended system which uses atomic scale modeling, involve extensive computations which requires supercomputer capability to achieve results in reasonable times. A high speed processor card has been developed which provides supercomputer capability on a PC platform and allows the MMR characterization systems to be implemented in a small package at low cost. The heart of the new computer technology is a parallel information processor board which can reside in any 80x86 PC, and is capable of 60 million floating point operations per second (MFLOPS). The PIP board also has 512 kilowords (32 bits/word) of high speed random access static memory on board and is expandable to 20 megawords. Up to eight PIP boards can be installed in a single PC to achieve 480 MFLOPS computing power, approximately one half the peak computation rate of a CRAY I. A user-friendly interface for this computing system is being developed which allows programming in Fortran, Basic, C, and a high level tool for integrating various software modules with flow charts. This PIP technology has been licensed to Scientific Data Systems Inc., which expects to begin commercial distribution in late 1992.

Each of the procedures which comprise the system shown in Figure 5 has been completed. The procedures will be integrated on a 386 PC with a single PIP card. We expect to incorporate an atomic scale model with the MMR characterization system within a year. The extended system using atomic scale modeling will require additional PIP cards.

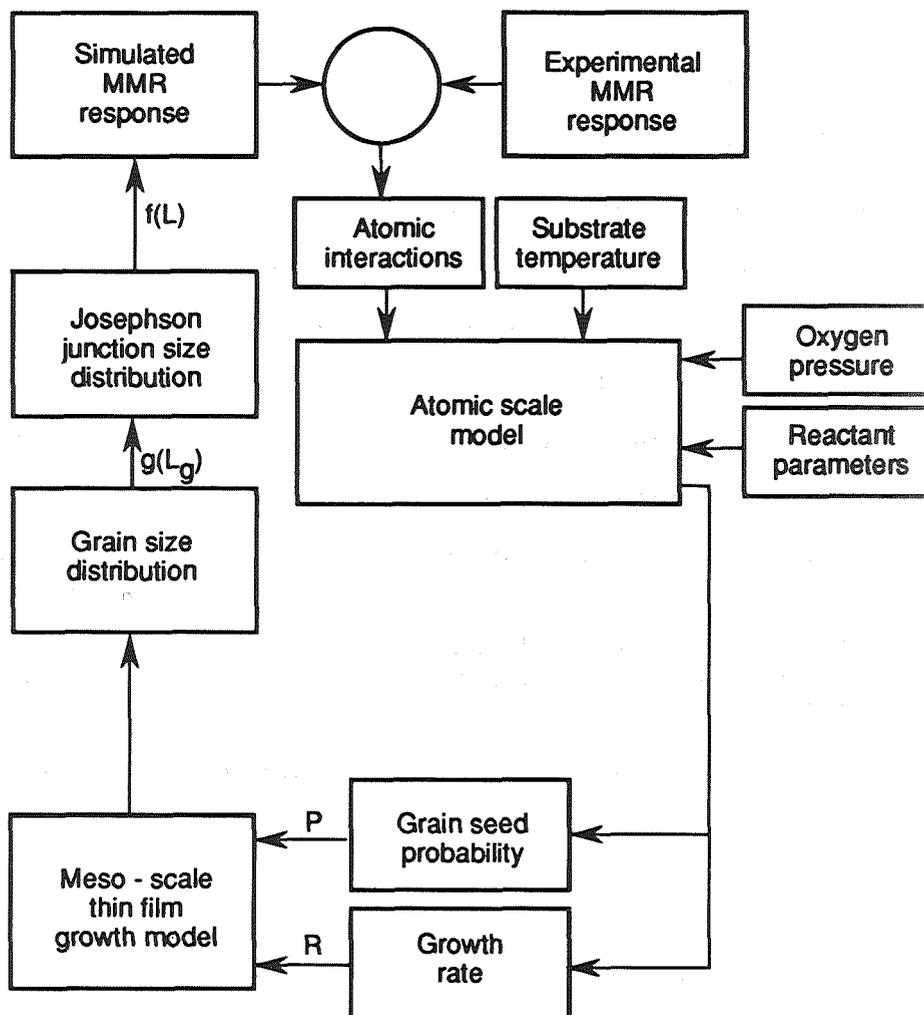


Figure 6. Schematic diagram of an MMR characterization system which uses atomic scale modeling.

CONCLUSION

The superconductor characterization system described here is noteworthy because of the extent of the information inferred from a relatively simple physical measurement. In its early usage, the MMR method provided evidence for the presence of weak links in superconductor samples in both granular samples and single crystals. At this stage, it was and remains a unique capability. The present stage of development in which the grain size distribution is inferred and is related to process parameters is due to the use of numerical simulations of the measured response function. This type of enhancement of physical measurements is being rendered practical by the emergence of powerful, low cost computer technology such as that embodied by the PIP card. There is little doubt that the combined use of measurements and large scale simulation in materials characterization and process control will become common as use of the new computer technology becomes widespread.

REFERENCES

1. B. F. Kim, J. Bohandy, K. Moorjani, F. J. Adrian, *J. Appl. Phys.* **63**, 2029 (1988).
2. J. Bohandy, T. E. Phillips, F. J. Adrian, K. Moorjani, B. F. Kim, *Mod. Phys. Lett. B* **3**, 933 (1989).
3. B. F. Kim, J. Bohandy, K. Moorjani, F. J. Adrian, . U. S. Patent #4,851,762, 1989.
4. J. Bohandy, B. F. Kim, T. E. Phillips, F. J. Adrian, K. Moorjani, U. S. Patent #4,904,929, 1990.
5. J. Bohandy, B. F. Kim, T. E. Phillips, F. J. Adrian, K. Moorjani, U. S. Patent #5,059,891, 1991.
6. Q. E. Dolecek, U. S. Patent #4,720,780, 1988.
7. Q. E. Dolecek, U. S. Patent #4,922,418, 1990.
8. Q. E. Dolecek, U. S. Patent #4,974,188, 1990.
9. See, for example, A. C. Rose-Innes and E. H. Rhodenck, "Introduction to Superconductivity," Pergamon Press, pp. 160-162 (1969).
10. B. F. Kim, K. Moorjani, F. J. Adrian, J. Bohandy, in *Magnetic Susceptibility of Superconductors and Other Spin Systems*, R. A. Hein, T. L. Francavilla, D. A. Liebenberg, Eds. (Plenum Press, New York, 1992) pp. 531-552.
11. B. F. Kim, J. Bohandy, F. J. Adrian, T. E. Phillips, and K. Moorjani, *Physica C* **161**, 76 (1989).

**PRODUCTION OF ULTRAFINE, HIGH-PURITY CERAMIC POWDERS
USING THE U.S. BUREAU OF MINES DEVELOPED TURBOMILL**

Jesse L. Hoyer
U. S. Bureau of Mines
Tuscaloosa Research Center
Tuscaloosa, AL 35486
(205) 759-9439

524-27
150494
P-10

ABSTRACT

Turbomilling, an innovative grinding technology developed by the U. S. Bureau of Mines in the early 1960's for delaminating filler-grade kaolinitic clays, has been expanded into the areas of particle size reduction, material mixing and process reaction kinetics. The turbomill, originally called an attrition grinder, has been used for particle size reduction of many minerals, including natural and synthetic mica, pyrophyllite, talc and marble. In recent years, an all-polymer version of the turbomill has been used to produce ultrafine, high-purity advanced ceramic powders such as SiC, Si₃N₄, TiB₂, and ZrO₂. In addition to particle size reduction, the turbomill has been used to produce intimate mixtures of high surface area powders and whiskers. Raw materials, TiN, AlN, and Al₂O₃, used to produce a titanium nitride/aluminum oxynitride (TiN/AlON) composite were mixed in the turbomill, resulting in strength increases over samples prepared by dry ball milling. Using the turbomill as a leach vessel, it was found that 90.4 pct of the copper was extracted from the chalcopryrite during a 4-hour leach test in ferric sulfate versus conventional processing which involves either roasting the ore for Cu recovery or leaching the ore for several days.

INTRODUCTION

The Bureau of Mines turbomill has evolved over a period of 30 years. The original patent granted to the Bureau was for a process to produce paper coating grade clays from lower grade kaolins [1]. Following this work, the Bureau conducted numerous other studies on grinding of industrial minerals. A Bureau of Mines Bulletin highlighting the kaolin research and information on the grinding of industrial minerals was published in 1980 [2]. This Bulletin also describes several commercial applications of the turbomill grinding technology.

In the early 1980's, research to determine the feasibility of using the turbomill to produce ultrafine, high-purity powders for advanced ceramics was undertaken [3]. In his study, Wittmer found that preparing SiC in the original steel mill yielded powders with iron levels above the acceptable range. A variety of polymers was tested as mill construction materials with ultra-high molecular weight (UHMW) polyethylene exhibiting the best wear resistance. Use of the all-polymer mill produced α -SiC powders of higher purity. Wittmer also evaluated the use of autogenous milling, in which the milling medium and the material to be milled are of the same or similar composition. Other ceramic materials, such as Si₃N₄ and ZrO₂, were milled with favorable results.

The turbomill also has been used to produce intimate homogeneous mixtures of high surface area powders and whiskers. SiC whiskers have been dispersed in alumina and silicon nitride powders [4]. Mixing the raw materials for preparation of a TiN/AlON composite in the turbomill resulted in strength increases over samples prepared by dry ball milling the components.

Rice, Cobble, and Brooks reported the use of the turbomill as a leaching vessel. Chalcopryrite was ground and leached simultaneously with ferric sulfate [5]. They found that 90.4 pct of the copper was extracted from chalcopryrite during a 4-hour leach test. Conventional processing technology for the recovery of Cu from chalcopryrite involves either roasting the chalcopryrite ore or leaching the ore for extended periods of time.

THE TURBOMILL

The turbomill consists of three main parts: a rotor, composed of vertical bars fixed to upper and lower disks (the upper one attached to the drive shaft); a cagelike stator composed of vertical bars attached to rings at the top and bottom; and a cylindrical container. A frame holds the motor which is attached to the rotor drive shaft and the machine components. The turbomill has been scaled up to sizes as large as 50.8-cm-diam for use in the laboratory. Industry has used mills as large as 132-cm-diam for production of ultrafine powders. The mill in use at the Tuscaloosa Research Center (TURC) is a 12.7-cm-diam. mill with the rotor, stator, and container constructed of UHMW polyethylene (shown in figures 1 and 2).

The slurry to be milled is placed in the container and consists of a milling liquid (typically water, kerosene, or alcohol), the milling medium (spherical balls or coarse granular material), and the material to be milled. A dispersant and/or antifoaming agent is typically added to aid in milling efficiency. The rotor operates at speeds which range from 1,400 to 1,700 rpm. Run times needed to reach the desired particle size range from 30 min to 4 hours.

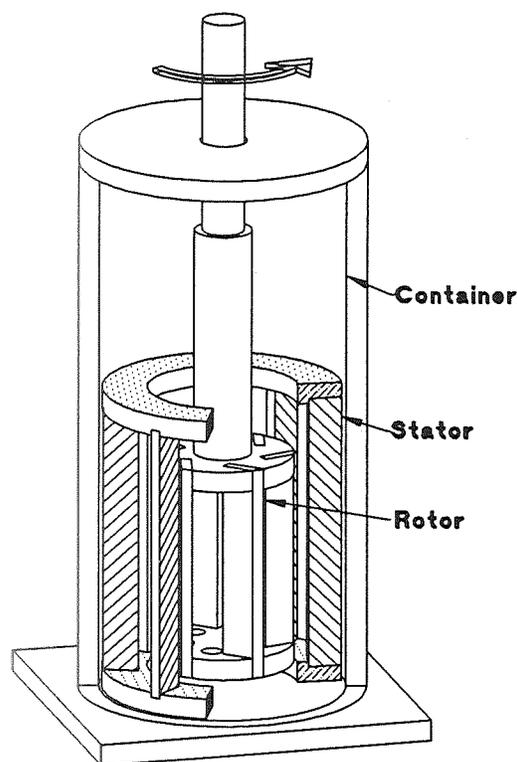


Figure 1. Bureau of Mines turbomill.

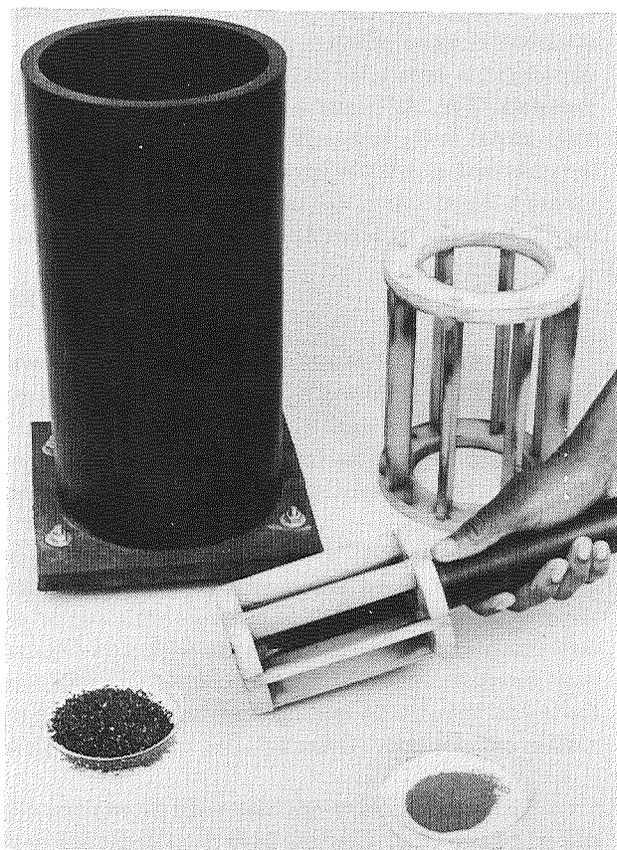


Figure 2. All-polymer container, stator, and rotor components of the turbomill.

GRINDING OF KAOLIN AND OTHER MINERALS

The efficiency of grinding kaolin using the turbomill was compared to the efficiency of other processes, including jar mills, colloid mills, and ultrasonics. In general, it was found that the turbomill produced more minus 2- μm equivalent spherical diameter (ESD) material than other methods in less time. It was also found

that the milling efficiency increases when a large mill (25.4-cm-diam) is used versus a smaller mill (12.7-cm-diam).

A series of tests to develop continuous grinding using the turbomill considered both open- and closed-circuit systems [6-7]. In open-circuit turbomilling, the overflow material was collected and sized. In closed-circuit milling, the overflow material was collected, sized, and incompletely ground material was returned to the mill as part of the feed slurry. It was found that closed-circuit grinding achieved a slightly greater rate of particle size reduction and consumed less energy than open-circuit grinding. A later investigation to determine the effect of operating variables on the grinding efficiency of kaolin showed the greatest influences related to the type, size and shape of grinding media; the ratio of medium weight to clay weight; peripheral rotor speed; clay slime pulp density; density of pulp dispersion and the angular arrangement of the rotor and stator. [8].

Following the grinding studies on kaolin, the Bureau determined the extent of particle size reduction obtained with other minerals [9]. A series of tests on mica showed that a sample with 11.1 pct minus 45- μ m material was ground with sand to a product 100 pct minus 45- μ m and 52.3 pct minus 2- μ m material. The electrical energy consumed when grinding the mica sample at a feed pulp density of 50 percent solids was 1206.6 MJ/mt of dry feed (304 kW•h/st) of dry feed. A second sample of mica with 8 pct minus 45- μ m material was ground to 100 pct minus 45- μ m with 56.8 pct minus 2- μ m material. Electrical consumption for the second sample was 1270 MJ/mt (320 kW•h/st) of dry feed. A large scale continuous open-circuit method for grinding mica was also developed. A feed with 9.2 pct minus 45- μ m mica was ground to 51.3 pct less than 45- μ m. Pyrophyllite, talc, marble, barite, and fluorite were also evaluated as part of this research.

BENEFICIATION APPLICATIONS

Other research focused on improved beneficiation of minerals including clay and olivine foundry sand. A study to improve the plasticity of a coarse kaolinitic clay for use as a bond clay in ceramic bodies was reported by Goode and Tyrrell [10]. An Alabama underclay was ground in the turbomill and compared to a standard ball clay used for whiteware bodies. The modulus of rupture of the fired body made using the ground underclay was increased by a factor of 1.7 over bodies made with the unground clay and was close to those made using the standard ball clay. The water absorption of the fired body made with turbomilled underclay was reduced by 35 pct when compared to the body made with the unground underclay.

A study to determine the feasibility of producing olivine for use as a foundry sand from dunite was reported in 1977 [11]. Dunite contains olivine and variable quantities of serpentine, talc, chlorite, actinolite, and vermiculite. The original products analyzed 2- to 3-pct loss on ignition (LOI) resulting from these impurities. After grinding using the turbomill, 200-mesh concentrates of the ground material met or exceeded the Steel Founder's Society of America specifications of a maximum 1.35 pct LOI for olivine aggregate and flour products. This beneficiation technique was patented by Davis in 1977 [12].

CERAMIC POWDERS

In the 1970's researchers at the Bureau's Tuscaloosa Research Center (TURC) investigated turbomilling of several ceramic oxide powders [13]. In general, the turbomilling process was more efficient than other processes such as ball milling. The increased particle size reduction yielded more reactive powders for sintering. For example, zirconia produced in the turbomill was sintered to 96 pct of theoretical density while commercial powders only sintered to 74 pct of theoretical [14].

Another research program was designed to demonstrate the catalytic properties of turbomilled silica [15]. Two crystalline forms of silica were ground and compared to two commercial silica catalysts used for n-hexane cracking and for dehydration of ethanol. The milled catalysts were more active than the commercial catalysts for n-hexane cracking; however, the ethanol dehydration activity was lower for the turbomilled silica. The differences were due to the different active sites required for dehydration.

Stanley, et al. describe a method for autogenous milling of SiC [16]. In these early tests, the mill was constructed of metal, resulting in iron contamination of the mill product. In 1980, a study to determine the feasibility of producing high purity α -SiC powders in the turbomill was undertaken [3]. Since iron is detrimental to the properties of sintered α -SiC, researchers redesigned the turbomill using a polymer as the construction material for the mill parts. Several polymers were tested and UHMW was selected for its wear resistance properties. Autogenous milling was used to eliminate contamination of the mill product by the milling media.

Autogenous turbomilling in the all-polymer turbomill successfully produced ultrafine α -SiC powders. Powders with Brunaur-Emmett-Teller (BET) surface areas of 30 to 35 m²/g were obtained after 3 to 6 hours of milling. These powders, with iron contents less than the starting material, were hot pressed with 1 pct boron (B) and 1 pct carbon (C) additions. Theoretical density of >99 pct was achieved and properties comparable to those of commercially available α -SiC were obtained.

Based on these promising results, experiments to optimize the milling parameters for α -SiC were conducted. The starting material for the tests had an average diameter of 100 μ m. Ninety-seven percent of the starting material was less than 150 μ m in diameter and contained no material below 30 μ m in diameter. The effects of dispersants, temperature, pH, and milling time on the particle size were investigated. The six dispersants used were: TSPP,* a sodium phosphate; Darvan No. 7,[†] a sodium salt; Marasperse N-22,[‡] a sodium lignosulfonate; Nopcosperse 44,[§] an ammonium salt; Aerosol OT,[¶] an anionic disodium sulfosuccinate; and Norlig NH,^{**} an ammonium lignosulfonate.

Table 1 lists average particle diameters, determined by laser beam diffraction, and ESD calculated from BET measurements of the α -SiC ground in water using the different dispersants. The average particle size produced using no dispersant, Marasperse N-22, Norlig NH, Aerosol OT, TSPP or Darvan No. 7 decreased gradually during 4 hours. Particle size did not change from 1 h to 4 h when Nopcosperse 44 was used as the dispersant; however, the ESD of the powder decreased. This indicates the formation of agglomerates during milling which was confirmed using scanning electron microscopy (SEM). Particle size distribution data, listed in table 2, for Marasperse N-22 indicates that the amount of material less than 1- μ m was 55 pct after 1 hour of milling; the amount does not increase significantly during the next 3 hours.

Contamination of the α -SiC, determined by spectrographic analysis, was negligible. The iron content decreased with milling time. During milling, the surface of the SiC particles is scrubbed resulting in removal of the iron, which is a surface contaminant. As the iron is removed from the surface of the powders, it goes into solution. The decrease in iron content was advantageous because SiC must often be leached after grinding to reduce the iron content.

* Fisher Scientific, Fair Lawn, NJ. Reference to a specific product does not imply endorsement by the Bureau of Mines.

[†]R. T. Vanderbilt Co., Norwalk, CT.

[‡] Reed Lignin, Inc., Atlanta, GA.

[§] Diamond Shamrock, Inc., Morristown, NJ.

[¶] American Cyanamid, New York, NY.

^{**} Reed Lignin, Inc., Atlanta, GA.

Table 1. Effect of dispersant on particle size of α -SiC as determined by laser beam diffraction (average diameter) or by surface area measurement (ESD) after milling for 1 hour and 4 hours.

Dispersant	Average diameter, μm		ESD, μm	
	1 h ¹	4 h	1 h	4 h
None	2.13	1.81	0.21	0.13
TSP	1.16	0.98	0.27	0.16
Darvan No. 7	1.54	0.99	0.22	0.16
Marasperse N-22	1.31	0.92	0.24	0.17
Nopcosperse 44	1.35	1.42	0.22	0.18
Aerosol OT	2.47	1.00	0.32	0.20
Norlig NH	3.22	1.92	1.05	0.28

¹Milling time

Note.--Average diameter of starting material - 100 μm

Table 2. Effect of Marasperse N-22 dispersant on α -SiC particle size

Particle diameter, μm	Percent ¹ less than indicated diameter			
	1 h ²	2 h	3 h	4 h
2.21	84	90	96	96
1.30	71	80	86	86
0.80	51	59	65	65
0.55	30	34	38	38
0.39	12	12	13	13
0.30	3	3	3	3
0.20	1	1	1	1

¹Percentages are for minus 325-mesh fraction of milled material.

²Milling time.

Note.--Average diameter of starting material - 100 μm

When Marasperse N-22 was used as the dispersant, the temperature of milling also played a major role in milling efficiency, as shown in table 3. As the temperature increased, the resulting average particle size did not change significantly. However, the grinding efficiency, as indicated by the amount of minus 325-mesh and submicrometer material, increased. The poor correlation between ESD values and the laser-measured particle diameters results from the fact that the laser technique measures agglomerates while BET techniques measure individual particles. Changes in pH also did not significantly affect the particle size (table 4); however, the efficiency of grinding was increased when the pH was slightly basic. The dispersing effect of Marasperse is affected by pH according to its manufacturers, with best results between pH 7 and pH 10. This increased dispersion effect results in greater particle/particle contact and increased grinding efficiency. The combined

effects of Marasperse N-22, pH 9.5 and 50° C resulted in 80 pct <1- μ m material after 4 hours of milling. The energy required was 2088.9 MJ/mt of feed (526.4 kW•h/st of feed).

Table 3. Effect of temperature on 4 hour milling of α -SiC in water with Marasperse N-22

Milling temperature, ° C	ESD, μ m	Average diameter, μ m	Percent less than		Contaminants, pct	
			325 mesh	1 μ m	Fe	Na
25	0.18	0.87	60.6	43.8	0.18	0.13
50	0.19	0.76	78.0	57.5	0.05	0.14
70	0.15	0.73	82.1	63.1	0.00	0.19

Note.--Average diameter of starting material - 100 μ m

Table 4. Effect of pH on 4 hour milling of α -SiC in water with Marasperse N-22

Milling temperature, ° C	pH	ESD, μ m	Average diameter, μ m	Percent less than	
				325 mesh	1 μ m
25	3.6	0.22	0.90	63.0	49.9
25	6.5	0.18	0.87	60.6	43.8
25	9.5	0.20	0.82	86.5	67.9
50	9.5	ND	1.14	92.6	79.7

Note.--Average diameter of starting material - 100 μ m

Several other materials used to produce advanced ceramics were also evaluated following turbomilling. Tables 5 and 6 list results for silica, two alumina materials, zirconia (ZrO_2), Al_2O_3 -partially stabilized zirconia (APSZ), CaO-partially-stabilized-zirconia (CPSZ), silicon nitride (Si_3N_4), and TiB_2 . These materials were ground in water using autogenous turbomilling. The average diameter and the amount of submicrometer material follow similar trends except for the larger Al_2O_3 material and the TiB_2 . The particle size of the Al_2O_3 decreased rapidly and then increased and the percentage of submicrometer material decreased. The increase in particle size and decrease in the amount of submicrometer powder indicates agglomeration to form larger particles or breakdown of the milling media. SEM analysis confirmed the presence of agglomeration and a sieve analysis of the mill product showed a 15 pct reduction in the amount of coarse milling media. Grinding of TiB_2 did not produce a large amount of submicrometer material. Use of other dispersants, different milling media or other milling fluids might increase the efficiency for milling of TiB_2 .

Table 5. Particle size of oxide ceramics

Material	Average diameter, μm				Percent less than 1- μm			
	0 h ¹	.5 h ²	1 h	4 h	0 h	.5 h	1 h	4 h
SiO ₂	16.3	ND	3.44	2.85	3	ND	14	28
Al ₂ O ₃ -1	102	ND	2.69	4.32	0	ND	31	4
Al ₂ O ₃ -2	34.2	1.96	ND	ND	0	24	ND	ND
ZrO ₂	20	ND	3.09	2.73	0	ND	17	23
APSZ	21.1	5.22	3.11	ND	0	16	23	ND
CPSZ	14	6.68	3.48	ND	4	10	17	ND

¹As-received particle size

²Milling time

ND-Not determined

The GTE Si₃N₄, listed in table 6, contained some whiskers. After 4 hours of milling, the whiskers were still present. The material was milled an additional 2 hours to determine if the whiskers would break down. After 6 hours, the powder no longer contained whiskers. Contamination of the powders was negligible as in the SiC studies.

Table 6. Particle size of nonoxide ceramics

Material	Average Diameter, μm				Percent less than 1- μm			
	0 h ¹	2 h ²	4 h	6 h	0 h	2 h	4 h	6 h
Si ₃ N ₄ [*]	55	2.09	2.05	ND	0	54	51	ND
Si ₃ N ₄ [‡]	55	2.72	2.51	2.53	0	23	35	37
TiB ₂ [§]	55	5.52	3.99	ND	0	3	4	ND

¹As-received particle size

²Milling time

^{*}Ube Chemical Co., Ube, Japan.

[‡]GTE, Towanda, PA.

[§]Sohio Chemicals and Industrial Products, Co., Niagara Falls, NY.

ND-Not determined

ALTERNATIVE PROCESSING OF CERAMIC RAW MATERIALS

The turbomill has also been used to improve the dispersion of sintering aids in high-surface-area powders [4] and to mix high surface area powders for a ceramic matrix composite [17]. Boron and carbon were added to α -SiC during turbomilling. These powders were hot pressed and physical properties were compared to B- and C-doped SiC materials prepared with powders processed using traditional methods. The addition of B and C during turbomilling resulted in enhanced properties, likely due to the breakdown of agglomerates of carbon which serve as flaw origins in conventionally processed powders. The raw materials, TiN, AlN and Al₂O₃ powders, for a ceramic matrix composite were mixed in the turbomill. Comparison of hot pressed composites formed from ball milled and turbomilled powders showed that turbomilling is the preferred method for preparing the powders for the composite.

Many advanced ceramic materials have whiskers added as reinforcements to improve the fracture toughness of the matrix material. A uniform dispersion of the whiskers is essential to the production of a reliable composite. Wittmer [4] conducted studies sponsored by Martin Marietta Energy Systems-ORNL to determine the feasibility of using the turbomill to produce a SiC whisker reinforced alumina composite. Well dispersed mixtures with no whisker agglomerates were prepared in 30 min in the turbomill using partially stabilized zirconia beads as the milling media. No degradation of the whiskers was visible in the short milling time. The strength and fracture toughness of these composites exceeded the properties measured on composites prepared by conventional methods by 30 pct. Wittmer also demonstrated that the slurries prepared by turbomilling could be pressure cast to form green bodies which, when fired, had improved properties over those processed conventionally.

MICROGRINDING OF COAL

The turbomill was used to produce microfine coal which could be used as a substitute for oil in firing steam boilers or as an addition to diesel fuel [18]. Three coals, each from a different seam, were ground in water using Ottawa sand, steel shot or coarse grain coal as the grinding media. Grinding with steel shot was the most efficient technique. A coal from the Pittsburgh seam (82 pct minus 75- μm) was reduced to 57 pct minus 2- μm in 15 min with energy requirements of 496.8 MJ/mt (138 kW•h/mt).

ALTERNATIVE LEACHING PROCESS FOR CHALCOPYRITE

Copper is traditionally recovered from chalcopyrite ore concentrates by roasting which results in sulfur dioxide emissions or by leaching the chalcopyrite ore for extended periods of time. In its efforts to develop new mineral processing technologies to enable cost-effective compliance with environmental requirements, the Bureau conducted research to improve the kinetics of chalcopyrite dissolution by simultaneous grinding-leaching in ferric sulfate solution [5]. The researchers proposed that during grinding, a chalcopyrite surface would be exposed, leached, and removed. Grinding of the surfaces in the turbomill would increase the number of fresh surfaces exposed to the leaching solution, reducing diffusion barriers for the reaction. Using minus 20-plus 30-mesh Ottawa sand as the grinding media, they found that the rate of leaching increased with increasing rotor speed and increasing solids content (by volume). They also found that the energy required for leaching passes through a minimum level as the amount of solids is varied. The minimum energy level for the conditions studied was 4,576 MJ/mt (1,153 kW•h/st) of copper extracted. This level occurred at 400 rpm and 35 vol pct solids at a recovery of 80 pct.

TECHNOLOGY TRANSFER AND COMMERCIAL APPLICATIONS

The turbomill developed by the U.S. Bureau of Mines has been used to produce materials for many different applications. It has been used in minerals processing. It has also been used to produce high-purity powders for use as raw materials for advanced ceramics. At the Tuscaloosa Research Center of the Bureau of Mines, a 12.7-cm-diam laboratory mill is available for evaluating different materials. Since 1987, the materials listed in table 7 have been milled in the turbomill. The Bureau of Mines can enter into various types of agreements in order to transfer this technology into different areas of interest. A Memorandum of Agreement (MOA) can be used to conduct work with private, state and academic organizations. An MOA outlines the work to be done and the costs involved. A Memorandum of Understanding (MOU) is similar to an MOA except it is an agreement with another federal agency. The Bureau can also use a Cooperative Research and Development Agreement (CRDA) with any non-Federal party. A CRDA is usually used when the cooperative work involves proprietary information and/or there is the possibility of a patentable invention.

Several small mills have been constructed by researchers around the world for use in their labs, while others have expressed interest in purchasing production size mills. Because these mills are built from schematic diagrams, the Bureau does not know how many units are in use worldwide. Turbomilling technology has been applied successfully on a commercial scale in the paper coating and titania pigment preparation industries. Georgia Kaolin Co. has designed and is operating 101.6-cm-diam production units, each capable of treating

1.81 mt of coarse kaolin/h (2 st of coarse kaolin/h). Kerr-McGee uses 132.1-cm-diam units with a combined capacity of 45,350 mt/yr (50,000 st/yr) to improve the particle size uniformity of TiO₂ for paint pigment.

Table 7. Materials evaluated in the Bureau of Mines turbomill

Material	Description/Application
Kaolin	Low grade, clay product. Desire to improve paper coating quality.
Alumina	Electrical applications
Graphite	Pencil Lead
Rutile waste	Desire to remove calcium carbonate from surface to allow TiO ₂ recovery
Hematite	Paint pigment
Y ₂ O ₃ stabilized ZrO ₂	Skull melted
Al ₂ O ₃ /ZrO ₂	Skull melted
Silica	Potential raw material for optical fiber or computer chip industry
B"-alumina	Electrolyte for use in fabrication of Na-S batteries
Indium oxide	Alloy manufacturing
Glass powders	Dielectric applications
Aluminum titanate	Paint pigment
Petroleum Coke	Residue from oil refinery. Recovery for use as particulate addition to rubber
Silica	Desire to produce silica flour from plus 100-mesh material
Silicon Nitride with MgO and Y ₂ O ₃	Obtain intimate mixture of MgO and/or Y ₂ O ₃ in Si ₃ N ₄
Plastic film and chopped plastic bottles	Desire to scrub surface of plastics to remove adhesive and printing prior to recycling.
Glass frit	Reduce plus 200-mesh material with little contamination to <2-μm
Graphite fibers	Desire to reduce length from 0.25 in to 25-30 μm with no degradation of fiber properties

BIBLIOGRAPHY

1. Feld, I. L. and B. H. Clemmons. Process for Wet Grinding Solids to Extreme Fineness. U.S. Pat. 3,075,710, Jan. 29, 1960.

2. Stanczyk, M. H. and I. L. Feld. Comminution by the Attrition Grinding Process. BuMines B 670, 1980, 43 pp.
3. Wittmer, D. E. Use of Bureau of Mines Turbomill to Produce High-Purity Ultrafine Nonoxide Ceramic Powders. BuMines RI 8854, 1984, 12 pp.
4. _____. Alternative Processing Through Turbomilling. Am. Cer. Soc. Bull., v. 67, No. 10, 1988, pp. 1670-1672.
5. Rice, D. A., J. R. Cobble, and D. R. Brooks. Effects of Turbomilling Parameters on the Simultaneous Grinding and Ferric Sulfate Leaching of Chalcopyrite. BuMines RI 9351, 1991, 9 pp.
6. Stanczyk, M. H. and I. L. Feld. Continuous Attrition Grinding of Coarse Kaolin, Part 1. Open Circuit Tests. BuMines RI 6327, 1963, 14 pp.
7. _____. Continuous Attrition Grinding of Coarse Kaolin, Part 2. Closed-Circuit Test. BuMines RI 6694, 1965, 13 pp.
8. _____. Investigation of Operating Variables in the Attrition Grinding Process. BuMines RI 7168, 1968, 28 pp.
9. _____. Ultrafine Grinding of Several Industrial Minerals by the Attrition Grinding Process. BuMines RI 7641, 1972, 25 pp.
10. Goode, A. H., and M. E. Tyrrell. Beneficiation of Alabama Clays. BuMines RI 8071, 1975, 7 pp.
11. Lamont, W. E., G. V. Sullivan, E. G. Davis, and S. D. Sanders. Olivine Foundry Sand from North Carolina Dunite by Differential Grinding. SME, AIME Preprint 77-H-369, 1977, 22 pp.
12. Davis, E. G. (assigned to U.S. Department of the Interior). Beneficiation of Olivine Foundry Sand by Differential Attrition Grinding. U.S. Pat. 4,039,625, Aug. 2, 1977.
13. Stanley, D. A., L. Y. Sadler III, and D. R. Brooks. Size Reduction of Ceramic Powders By Attrition Milling. Proc. of IITRI Conf. in Part. Tech., Chicago, IL, Aug. 20-24, 1973, pp. 280-285.
14. Stanley, D. A., L. Y. Sadler III, D. R. Brooks, and M. A. Schwartz. Attrition Milling of Ceramic Oxides. Am. Ceram. Soc. Bull., v. 53, No. 11, 1974, pp. 813-815.
15. Hatcher, W. J., Jr., and L. Y. Sadler, III. Catalytic Properties of Attrition Ground Silica. J. of Catalysis, v. 38, 1975, pp. 73-39.
16. Stanley, D. A., L. Y. Sadler III, D. R. Brooks, and M. A. Schwartz. Production of Submicron Silicon Carbide Powders by Attrition Milling. Pres. at Int. Symp. on Fine Particles, Ann. Meet. of Electrochem. Soc., Boston, MA, Oct. 7-12, 1973, 14 pp.
17. Hoyer, J. L., J. P. Bennett, and K. J. Liles. Properties of TiAlON/Spinel Ceramic Composites. Ceram. Eng. Sci. Proc., v. 11, No. 9-10, 1990, pp. 1423-1439.
18. Davis, E. G. Fine Grinding of Coal by the Turbomilling Process. BuMines RI 9070, 1987, 8 pp.

MULLITE WHISKERS AND MULLITE-WHISKER FELT

Inna G. Talmy and Deborah A. Haught
 Naval Surface Warfare Center, Dahlgren Division
 Silver Spring, MD 20903-5000

525-24
 150495
 P-9

ABSTRACT

Naval Surface Warfare Center has developed processes for the preparation of mullite ($3\text{Al}_2\text{O}_3 \cdot 2\text{SiO}_2$) whiskers and mullite-whisker felt. Three patents on the technology were issued in 1990. The processes are based on chemical reactions between AlF_3 , Al_2O_3 , and SiO_2 . The felt is formed in-situ during processing of shaped powdered precursors. It consists of randomly oriented whiskers which are mutually intergrown forming a rigid structure. The microstructure and properties of the felt and size of the whiskers can be modified by varying the amount of Al_2O_3 in the starting mixture. Loose mullite whiskers can be used as a reinforcement for polymer-, metal-, and ceramic-matrix composites. The felt can be used as preforms for fabricating composite materials as well as for thermal insulation and high temperature chemically stable filters for liquids (melts) and gases.

INTRODUCTION

The current interest in the development of refractory-oxide fibers and whiskers is spurred by the increasing demand for high-temperature structural materials for use in an oxidizing atmosphere. As a reinforcing material, whiskers are preferable to fibers because they are single crystals and their properties are not affected by grain growth and grain boundary-induced creep at high temperatures. The low free energy and high modulus and strength of whiskers compared to polycrystalline materials make it possible to use whiskers to reinforce matrices of the same composition.

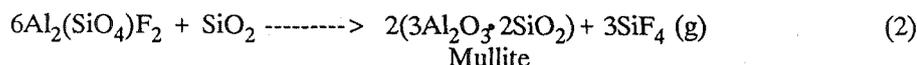
Currently, only non-oxide whiskers such as SiC and Si_3N_4 are commercially available. However, they cannot be used in an oxidizing atmosphere at temperatures above 1000°C . Additionally, SiC whiskers are electrically conductive and not suitable for use in dielectric ceramics.

Mullite ($3\text{Al}_2\text{O}_3 \cdot 2\text{SiO}_2$) is a promising candidate for whiskers because of its excellent chemical stability, low thermal expansion, and good high-temperature strength and creep resistance. Since mullite has a relatively low dielectric constant and loss tangent, the whiskers may be suitable for toughening materials for dielectric applications (such as radomes).

Both loose whiskers (or fibers) and fibrous preforms can be used as reinforcement in the fabrication of composites. The use of fibrous preforms instead of loose whiskers (or fibers) excluded problems of their deagglomeration and uniform distribution in matrix materials. Matrices can be infiltrated into the preforms by vacuum impregnation, sol-gel processing, chemical vapor infiltration and melt infiltration techniques.

PREPARATION OF MULLITE WHISKERS AND FELT

The NAVSWC method for preparation of both mullite whiskers and felt is based on the following two reactions¹⁻⁵:



In reaction (1), AlF_3 and SiO_2 are heated at 700 - 900° C to form topaz as an intermediate product and in reaction (2), topaz is thermally decomposed at 1250 - 1400° C to yield mullite. The whiskers are grown as a result of a vapor-phase chemical reaction. Heating in a SiF_4 atmosphere is necessary for the preparation of high-quality whisker product. Up to 75% Al_2O_3 can be substituted for AlF_3 in the starting materials. In this case, the mullite whiskers form as a result of chemical interactions involving all of the starting components, as well as the gaseous SiF_4 generated in reactions (1) and (2) - assuming that corresponding portions of AlF_3 required for the process are produced by a reaction between SiF_4 and Al_2O_3 .

The whisker preparation process involves mixing the raw powders and firing the loose powders in closed SiF_4 -containing system. Mullite whiskers prepared at 1250° C in optimum SiF_4 atmosphere are shown in Figure 1. As produced, they form loose aggregates which are easily separated. The whiskers of rectangular cross section have a narrow size distribution with average aspect ratio of 30.

The felt preparation process involves mixing and shaping (by any conventional method) the powdered precursors (AlF_3 , SiO_2 and optional Al_2O_3), followed by firing at 1250 - 1400° C. A significant advantage of this process is that there are no whiskers in the precursor mixture. The whiskers are formed inside the product during firing. Because loose respirable whiskers are not handled at any step of the process, health hazards associated with handling whiskers are eliminated. The final felt product is about 80% porous with low dimensional changes (about 1% expansion) compared to the green shapes which is a valuable feature for the preparation of near-net shape composite preforms. Figure 2 shows the microstructure of the felt consisting of randomly oriented (in 3 dimensions) and uniformly distributed individual mullite whiskers and spherulites which are mutually intergrown or mechanically interlocked. With a bending strength (3 point) of about 3 MPa (427 psi), the felt is rigid and can be used in composite processing utilizing various impregnation techniques.

To decrease the amount of relatively expensive AlF_3 and accordingly decrease the amount of hazardous SiF_4 generated by the reactions, and possibly modify the size of the whiskers, Al_2O_3 (avg. particle size 0.05 μm) was substituted for AlF_3 in amounts up to 75 mole %. The whisker size gradually decreases with increasing Al_2O_3 content up to 75% Al_2O_3 for both loose whiskers and felt (Figure 3). In the felt, the number of spherulites also gradually decrease with increasing Al_2O_3 content. The linear expansion and porosity of the felt slightly decrease but the bending strength significantly increases with more than 25% Al_2O_3 substitution (Table 1.). The almost two-fold increase in strength with 75% Al_2O_3 substitution can be attributed to much smaller whiskers in the felt structure. The felt has pores of 1 to 30 μm depending on the $\text{AlF}_3/\text{Al}_2\text{O}_3$ ratio with very narrow size distribution (Figure 4).

Table 1. - Properties of Mullite-Whisker Felt With Substitutions of Al_2O_3 for AlF_3

Substituted Al_2O_3 (%)	Expansion (%)	Porosity (%)	Bending Strength (MPa/Psi)
0	1.0	78.5	3.1/441
25	1.0	78.1	2.9/412
50	0.5	77.4	5.4/768
75	0.0	76.3	7.0/995

APPLICATIONS OF MULLITE WHISKERS

Mullite whiskers have the potential to be used as reinforcement in a variety of polymer, ceramic and

metal matrix composites. Unfortunately, the full potential of these whiskers as a reinforcement has not been completely exploited. Mullite whiskers were successfully used as reinforcement for celsian ($\text{BaO} \cdot \text{Al}_2\text{O}_3 \cdot 2\text{SiO}_2$) and fused silica (SiO_2) composites⁶. The loose whiskers were uniformly mixed with the ceramic powders and then hot pressed to full density. The fracture toughness of fused silica was doubled by the addition of 20 vol% mullite whiskers with a cross section of 4-12 μm . For the celsian matrix, 20 vol% 12-26 μm whisker additions increased the toughness by 40%. Properties of the composites are shown in Table 2.

Table 2. Properties of Mullite Whisker Composites

Matrix	Mullite Whisker Loading (%)	Relative Density (%)	Fracture Toughness* ($\text{MPa}\sqrt{\text{m}}$)	Flexural Strength (MPa)
Silica	0	94.5	0.92	67
Silica	20% 4-12 μm	96.2	1.80	49
Celsian	0	99.1	1.56	115
Celsian	20% 12-26 μm	99.4	2.15	115

* Short Chevron Notch Beam Method

APPLICATIONS OF MULLITE-WHISKER FELT

The mullite-whisker felt can be used as a preform for the preparation of polymer-, metal-, and ceramic-matrix composites or by itself as thermal insulation or high temperature filters for liquids or gases. It can be infiltrated with a variety of matrices using various processing techniques such as vacuum impregnation, chemical vapor infiltration (CVI), and melt (glass or metal), sol-gel or pre-ceramic polymer infiltration to produce composites. Polymers with viscosities of up to 10,000 cp have been used to infiltrate the felt. Researchers at Ceramtec have demonstrated that mullite whisker composites can be fabricated using sol-gel techniques⁷. After 8 infiltration cycles, the density of the sample was 2.06 g/cm^3 with 17% porosity. The process is being optimized to yield composites with densities greater than 95% dense.

A forced flow thermal gradient chemical vapor infiltration technique developed at Oak Ridge National Laboratory was used to infiltrate the felt with a SiC matrix. Since the deposition parameters used for the infiltration experiment were optimized for a Nicalon preform, the felt was only infiltrated to about 85% density with variable densities through the thickness of the sample. The properties of the produced composites are given in Table 3.

Fully dense mullite-whisker felt composites were prepared by a vacuum and isostatic pressure-assisted infiltration technique using cast aluminum alloy A356⁸. The alloy was preheated to 680°C and the felt was preheated to 550°C . A 800 psi pressure of nitrogen was required to fully infiltrate the preform. Due to the lack of bonding between A356 Al and mullite, the mechanical properties of the composite did not reflect the reinforcing potential of the whiskers. There was little transfer of the load from the matrix to the mullite whiskers. It will be necessary to increase the affinity of the matrix towards the reinforcement. Many options exist in matrix selection; addition of a small percentage of Li or Mg to the matrix may be sufficient, or contemporary Al-Li alloys such as 2090, 8090, 2091 etc. may be useful in enhancing bonding at the whisker/matrix interface without degrading the strength of mullite.

Table 3. Properties of SiC Infiltrated Mullite-Whisker Felt

Specimen Position	Theoretical Density (%)	Flexural Strength (MPa)	Fracture Toughness* (MPa√m)
Top	84.6	79.6	2.18
Middle	86.9	97.2	3.15
Bottom	85.8	105.8	3.35

* Single Edge Notch Beam Method

The felt was also evaluated for possible use as a filter in coal gasification processing. Conversion of coal to a gaseous fuel requires filtering of hazardous impurities from final gases to meet environmental and turbine equipment requirements in advanced coal-fueled power generation systems. The criteria for successful use and operation of porous filters requires not only thermal, chemical, and mechanical stability of the material, but also long-term structural durability and high reliability. Preliminary results indicate that the mullite-whisker felt is a very promising candidate. It showed good chemical stability in the tested gas atmosphere and had very high collection efficiency.

SUMMARY

Processes for the preparation of mullite whiskers and rigid mullite-whisker felt have been developed and patented. The 80% porous felt consists of randomly oriented mullite whiskers mutually intergrown, forming a rigid structure. The felt is formed in-situ during processing shaped powder precursors; thereby eliminating the health hazards associated with handling loose whiskers. Both the mullite whiskers and mullite-whisker felt are promising new reinforcements for polymer-, metal- and ceramic-matrix composites. Additionally, the felt can be used by itself as thermal insulation or high temperature chemically stable filters for liquids (melts) and gases such as for coal gasification.

ACKNOWLEDGEMENTS

The authors would like to acknowledge the assistance of Dr. M. Norr in the SEM studies. We would also like to thank Rick Lowden of Oak Ridge National Laboratory for infiltrating the felt using his CVI process and characterizing the composites. This work was funded by the Office of Naval Technology through the Weapons and Spacecraft Materials Block Program headed by Dr. W.T. Messick.

REFERENCES

1. Talmy, I. and Haught, D., "Preparation of Mullite Whiskers", U.S. Patent No. 4911902, 27 March 1990.
2. Talmy, I. and Haught, D., "Preparation of Mullite Whiskers From AlF_3 , SiO_2 , and Al_2O_3 Powders," U.S. Patent No. 4910172, 20 March 1990.
3. Talmy, I. and Haught, D., "Rigid Mullite Whisker Felt and Method of Preparation", U.S. Patent No. 4948766, 14 August 1990.

4. I. Talmy and D. Haught, "Preparation of Mullite Whiskers", *Proceedings of Fiber-Tex 87 Conference*, NASA Conference Publication 3001, p. 69-78 (1987).
5. I. Talmy and D. Haught, "Preparation and Properties of Rigid Mullite-Whisker Felt", *Proceedings of the 12th Conference on Composite Materials and Structures*, NASA Conference Publication 3018, p. 1-11 (1988).
6. C.L. Conner, "Mullite Whisker Composite Fabrication", NAVSWC TR 90-344 (12 August 1990).
7. I.G. Talmy, D.A. Haught, and S. Limaye, "Porous Mullite Whisker Felt for Preshaped Composite Preform", *Materials and Processing Report*, Vol. 4, No. 7, 3 (1989).
8. D. Haught, I. Talmy, D. Divecha and S. Karmarkar, "Mullite Whisker Felt and Its Application in Composites", *Materials Science and Engineering*, A144, p. 207-214 (1991).

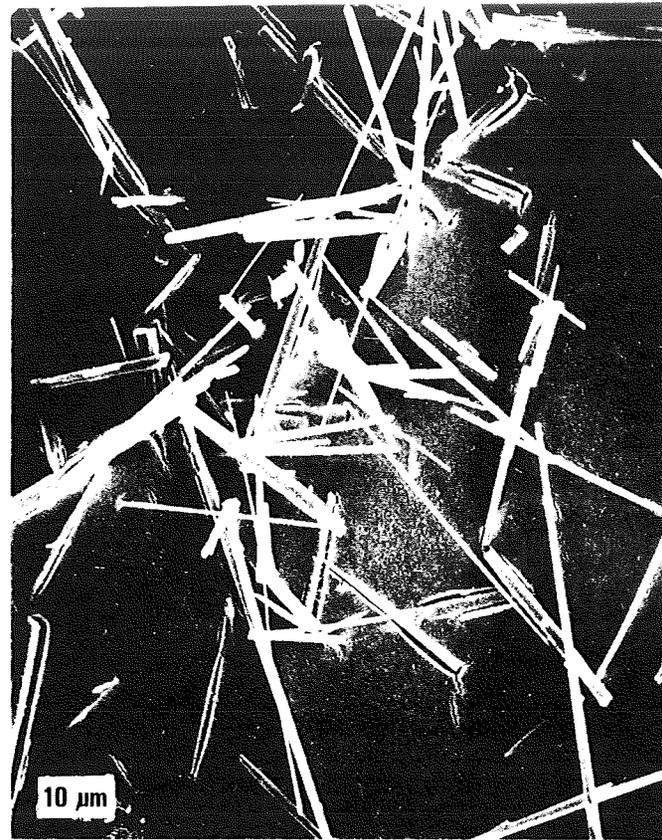


FIGURE 1. SCANNING ELECTRON MICROGRAPHS OF MULLITE WHISKERS



FIGURE 2. SCANNING ELECTRON MICROGRAPHS OF MULLITE WHISKER FELT

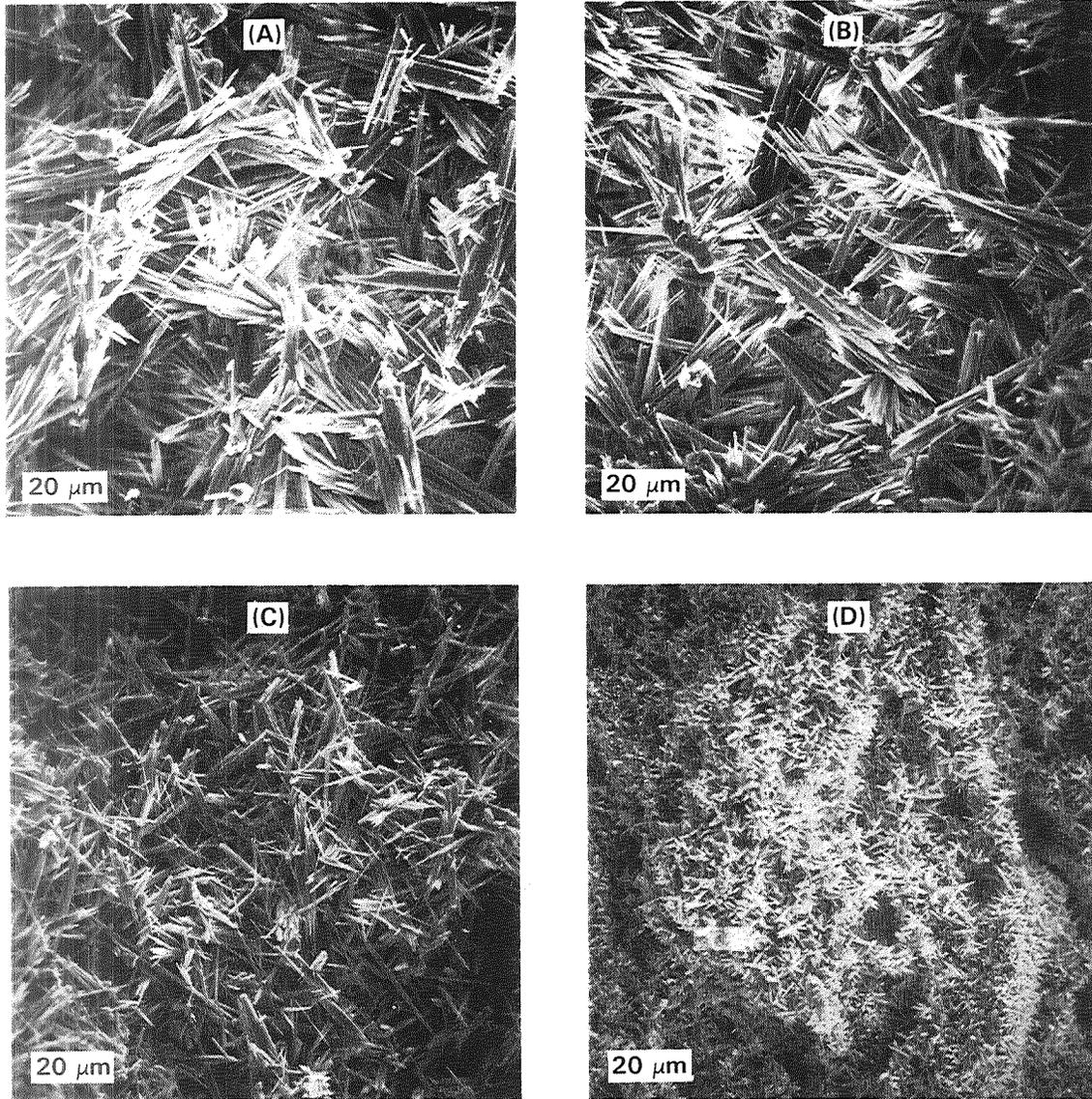


FIGURE 3. SCANNING ELECTRON MICROGRAPHS OF MULLITE WHISKER FELT WITH SUBSTITUTION OF (A) 0%, (B) 25%, (C) 50%, AND (D) 75% Al_2O_3 FOR AlF_3

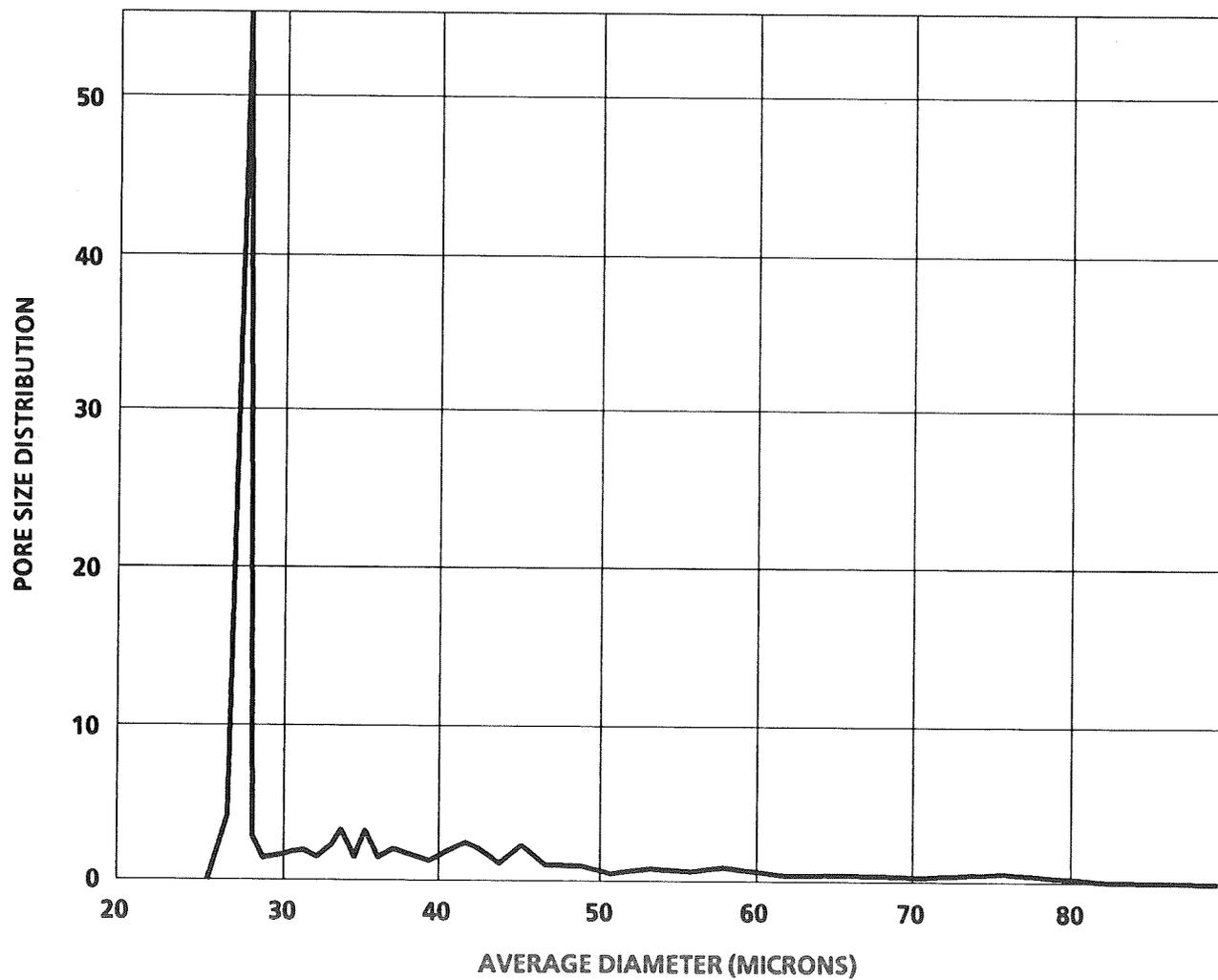


FIGURE 4. PORE SIZE DISTRIBUTION OF MULLITE WHISKER FELT PREPARED FROM AlF_3 AND SiO_2

526-24
150496
P-4

GRAPHITE/EPOXY COMPOSITE LAMINATES WITH
CO-CURED INTERLAMINAR DAMPING LAYERS

N 93-25587

J. Michael Pereira
NASA Lewis Research Center
21000 Brookpark Road, MS 49-8
Cleveland, Ohio 44135

ABSTRACT

Damped composite laminates were fabricated by co-curing viscoelastic damping film with graphite/epoxy prepreg plies. The dynamic response of the damped plates was measured using an impulse response technique and compared with the response of similar undamped laminates. Modal damping was computed from the frequency response data. Micrographs of the damped laminates showed that the damping layers retained their integrity during the fabrication process. The layers significantly increased the damping in the composite laminates. The use of the constrained viscoelastic film as an integral part of composite structures appears to be a feasible approach to passive vibration control. Composite plates manufactured with co-cured damping layers may have commercial applications in cases where light weight, strength and vibration and noise reduction are important considerations.

INTRODUCTION

Vibration control is an issue of prime importance in many structures subjected to external loads or internally moving components. Often, the control of vibration is of such critical importance that active vibration control systems are incorporated into the structure. While they may be effective, the use of active control systems can add significant complexity, weight and cost to a design. Whenever possible, it is usually more desirable to reduce vibration problems using passive methods such as dampers, additive damping material, or highly damped structural materials.

Composite materials have important potential benefits in the area of passive vibration control. Advanced composite materials can have significantly greater specific stiffness and damping than traditional structural materials such as aluminum [1, 2]. An important aspect of composites is that, like the stiffness properties, the damping capacity of composites can be controlled by the selection of materials and layup. Composite material systems can be engineered to control specific vibration problems while addressing other concerns such as stiffness, strength and toughness. Previous experiments have shown that the interlaminar fracture toughness of a graphite/epoxy laminate is increased by as much as a factor of 10 by curing an adhesive interply layer in the composite [3]. In addition, recent analytical studies have shown that the use of a constrained viscoelastic layer in a composite laminate can significantly increase modal damping [4].

The aim of this study was to experimentally investigate the use of a constrained viscoelastic damping material for passive vibration control of graphite/epoxy composite structures. The damping material was co-cured with the composite to form laminates with internal damping layers at various locations. The dynamic response of the composite plates was measured using an impulse response technique.

METHODS

Fabrication of Damped Laminates

Composite laminates with various layups were manufactured by co-curing graphite/epoxy composite in pre-preg form with a polymer damping material. The laminates consisted of T300/934 graphite/epoxy with a 60% fiber volume ratio combined with Scotchdamp ISD110 damping film (3M Corp., St. Paul, Minnesota). Two different laminates were fabricated with the damping film, a $[+45_2/-45_2/i/+45_2/-45_2]_{sym}$ laminate, and a $[+22.5_2/-22.5_2/i/+22.5_2/-22.5_2]_{sym}$ laminate, where i refers to the interlaminar damping layer. Two additional laminates were fabricated with the same layup, but excluding the damping layers, and one laminate with interlaminar damping layers was fabricated for microscopic examination. All laminates were cured in a press at 175°C (350°F) and 345 kPa (50 psi). After curing, the laminates were trimmed to produce 28cm x 28cm plate specimens.

Vibration Experiments

The panels were vibration tested in the free-free mode using an impulse response technique. The panels were supported from one corner by a very flexible rubber band, such that the fundamental mode of the panel/rubber band system had a frequency of less than 1 Hz. Impulse response tests were conducted by impacting the center of the plate with an instrumented hammer, and measuring the response with an accelerometer at one of four corner locations shown in Figure 1. The four locations were chosen to check the repeatability of measurements. The accelerometer had a mass of 1.0 gm, compared to a mass of approximately 270 gm for the plate specimens.

The recorded impulse response data consisted of the averaged force and acceleration history from 10 impacts. Data was recorded with a digital data acquisition system using an acquisition rate of 10^5 samples/sec and a record length of 16384 samples, for a total duration of .16384 sec. The data was transferred to a computer to calculate the frequency response functions and damping in the frequency range up to approximately 1500 Hz.

Damping Computation

For the purpose of computing damping, it was assumed that in the vicinity of a resonance the response was dominated by a single mode, which could be modeled as a single degree of freedom, viscously damped system. A modification of the standard Nyquist circle fit method [5] was developed to compute the specific damping capacity (SDC), ψ [6].

The effect of the damping layer on the plates was evaluated by comparing the SDC for damped and undamped plates as a function of frequency. The accuracy of the damping results was then evaluated by synthesizing the frequency response function using the computed modal constants, and comparing the result with the measured data.

RESULTS

There appeared to be no damage to the interlaminar damping layers as a result of the fabrication process. Micrographs of sections cut from one of the plates showed that the damping film retained its integrity, with no apparent delamination between the film and the adjacent composite plies (Figure 2). Rough areas visible in the central region of the damping layer shown in Figure 2 are believed to be a result of the micrograph specimen preparation process.

The synthesized frequency response functions were very close to the measured data (Figure 3), indicating that the assumed modal model was a good representation for the behavior of both the damped and undamped laminates.

The modal specific damping capacity for the undamped laminates was relatively constant over the frequency range examined (Figures 4 and 5). At the low end of the frequency range the SDC for the damped and undamped plates were comparable. The SDC for the damped plates increased with frequency, however, so that at the upper end of the frequency range examined the SDC for the damped laminates was 2 to 3 times that of the undamped laminates. The $\pm 22\frac{1}{2}$ laminates displayed some highly damped low frequency modes. It is not known whether this was a real phenomenon related to the material, or if it was an artifact due to rotational vibration of the plate/rubber band system. Apart from these low frequency modes, results for both layups were very similar, in both the damped and undamped laminates.

A qualitative measure of the effect of the damping layer can be obtained by comparing the acceleration time histories of the damped and undamped plates. The acceleration response decayed much faster with the damped plate than with the undamped plate (Figure 6).

DISCUSSION AND CONCLUSIONS

The approach of co-curing a damping film with composite pre-preg appears to be very effective as a passive vibration control measure. The damping layers significantly increase the modal damping in graphite/epoxy laminates which are already relatively highly damped. Although it required a number of trials to learn how to make the laminates, the processing is relatively straightforward. Micrographs indicate that the damping layers maintain their integrity and bond well to the adjacent composite plies.

While the constrained viscoelastic layers increase the structural damping in the laminates, further work should be done to determine how the layers influence the compressive strength, stiffness and impact resistance of the materials.

Determining the effect of a constrained viscoelastic layer on the dynamic response of composite structures is complex. The energy dissipation comes primarily through shear deformation of the layer. An important advantage of composite materials is that it is possible to tailor the material to optimize the in-plane shear and resulting damping. Methods of predicting the damping capacity of composites, and the damped dynamic response of structures are important if their benefits are to be realized.

REFERENCES

- [1] Crawley, E.F. and Van Schoor, M.C., 1987, "Material Damping in Aluminum and Metal-Matrix Composites," *Journal of Composite Materials*, Vol. 21, No. 6, pp. 553-568.
- [2] Bicos, A.S., 1989, "Damping of a Large Space Platform," *Proceedings, Damping 1989 Vol. II*, Report WRDC-TR-89-3116, Wright Research Dev. Ctr., pp. HBB-1 - HBB-9.
- [3] Browning, C.E. and Schwartz, H.S., 1986, "Delamination Resistant Composite Concepts," *Composite Materials: Testing and Design (Seventh Conference)*, ASTM STP 893, J.M. Whitney, ed., American Society for Testing and Materials, Philadelphia, pp. 256-265.
- [4] Saravanos, D.A. and Pereira, J.M., 1992, "The effects of Interply Damping Layers on the Damping Characteristics of Composite Plates," *ALAA Journal*, in press.
- [5] Ewins, D.J., 1984, *Modal Testing: Theory and Practice*, Research Studies Press Ltd., Herts, England.
- [6] Pereira, J.M., "Dynamic Response of Composite Plates with Interlaminar Damping Layers", *International Symposium on Vibroacoustic Characterization of Materials and Structures*, ASME Winter Annual Meeting, Anaheim, 1992.

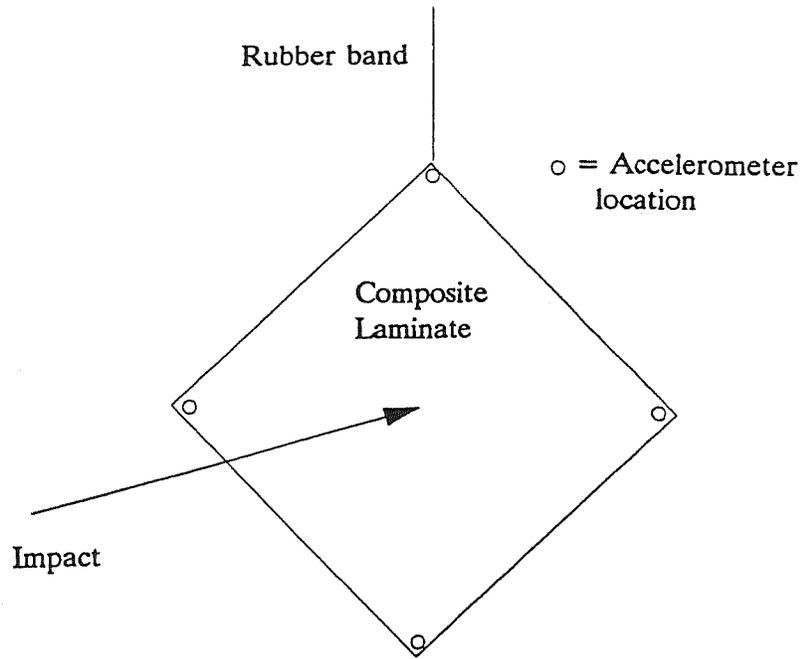


Figure 1: Schematic of the experimental configuration. Separate experiments were done with the accelerometer at each of the four locations shown.

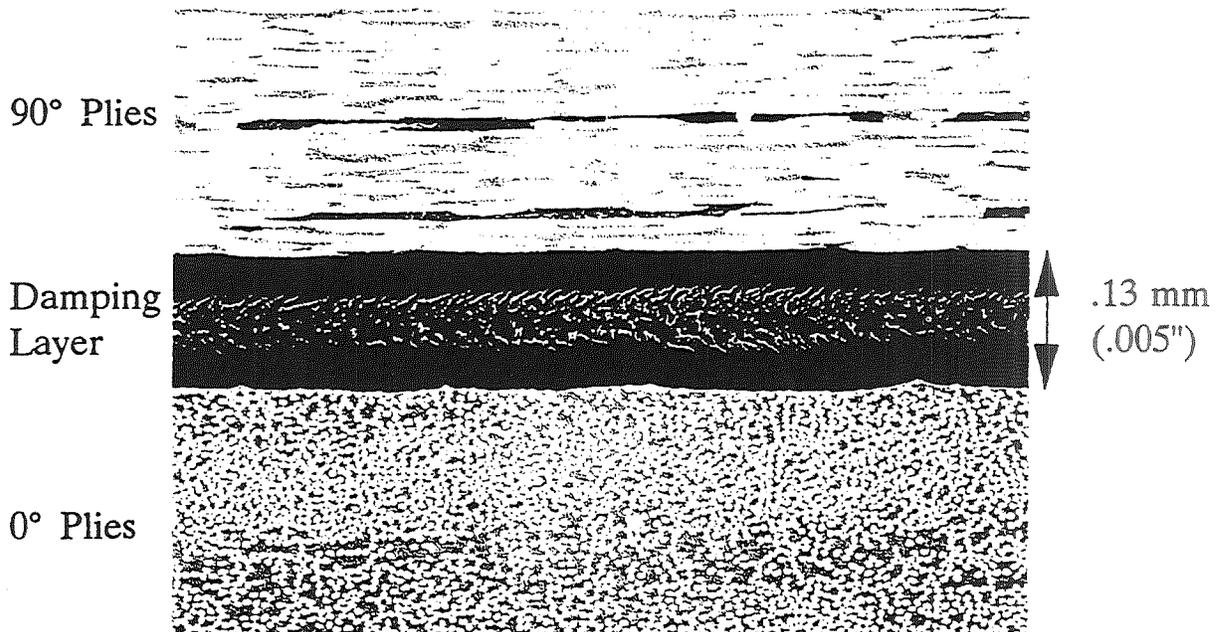


Figure 2: Micrograph of the damping layer between composite plies after the fabrication process. The damping layer appears to retain its integrity. The rough area at the center of the layer is believed to be a result of the micrograph specimen preparation process.

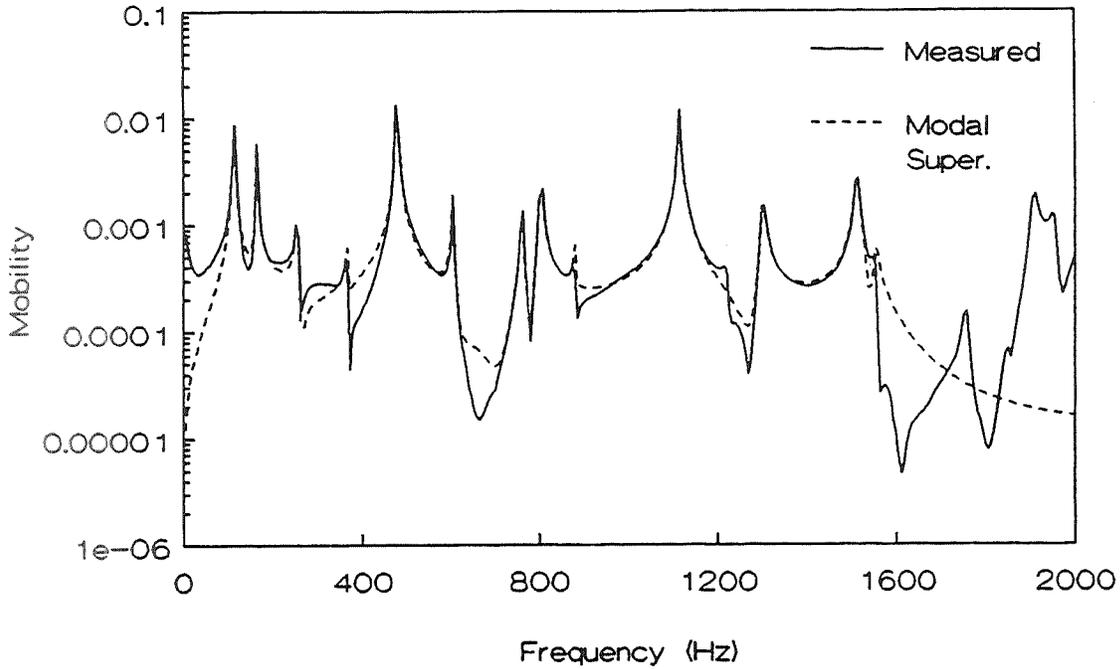


Figure 3: Comparison of the measured mobility function (velocity/force) with the predicted function using the modal superposition procedure for a ± 45 laminate without damping layers.

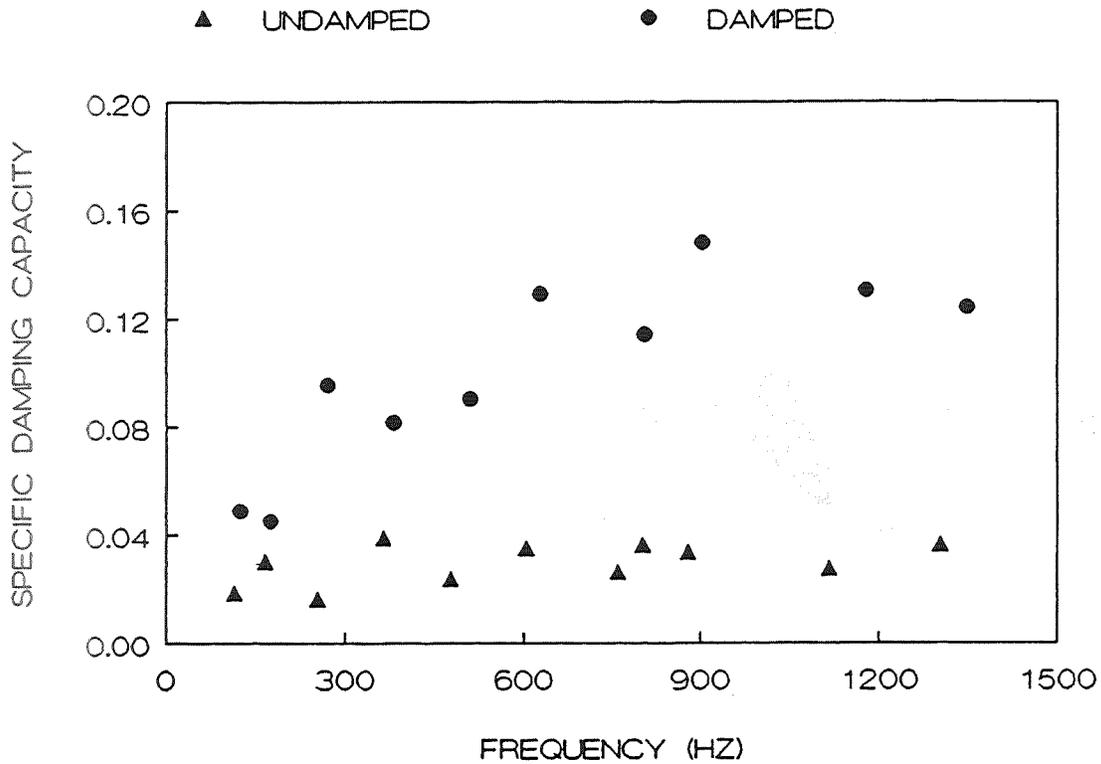


Figure 4: Modal specific damping capacity computed from the experimental data for the damped and undamped $[+45_2/-45_2/(i)/+45_2/-45_2]_{sym}$ laminates. The specific damping capacity for the damped laminate is significantly higher than that of the undamped laminate, and increases at higher frequencies.

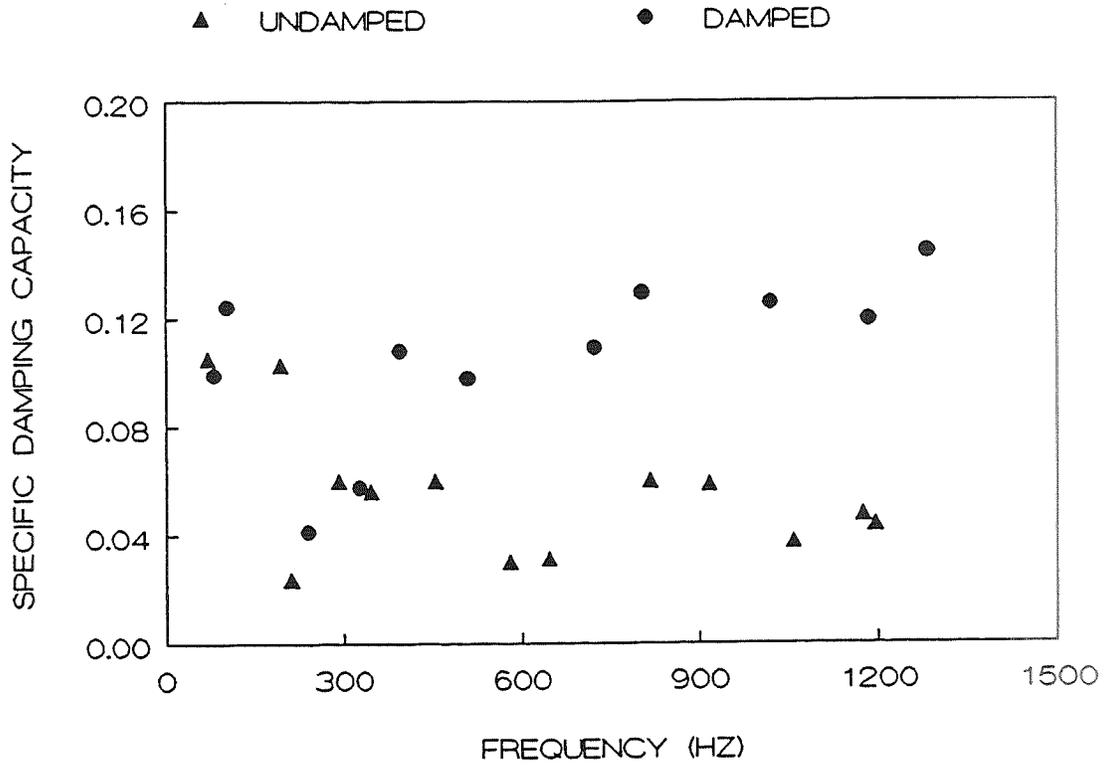


Figure 5: Modal specific damping capacity computed from the experimental data for the damped and undamped $[+22.5_2/-22.5_2/i/+22.5_2/-22.5_2]_s$ laminates. Apart from low frequency region, the specific damping capacity is significantly higher for the damped laminate.

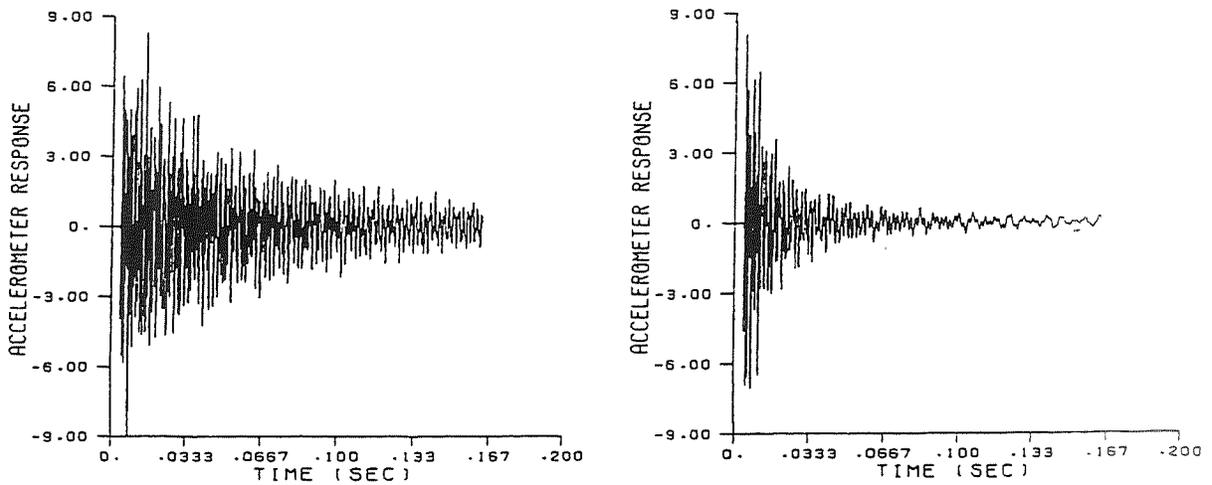


Figure 6: Accelerator time response for the undamped (left) and damped (right) $\pm 45^\circ$ laminates. There is a clear increase in the rate of decay for the damped laminate.



omit

**ARTIFICIAL INTELLIGENCE
PART 1**

PRECEDING PAGE BLANK NOT FILMED

257.

256 ~~INTENTIONALLY BLANK~~



EXPERT SYSTEM FOR UNIX SYSTEM
RELIABILITY AND AVAILABILITY ENHANCEMENT

Catherine Q. Xu, Ph.D.
Senior Staff Engineer
Aeronautical Radio, Inc.
Annapolis, MD 21401

ABSTRACT

Highly reliable and available systems are critical to the airline industry. However, most off-the-shelf computer operating systems and hardware do not have built-in fault tolerant mechanisms, the UNIX workstation is one example. In this research effort, ARINC has developed a rule-based Expert System (ES) to monitor, command, and control a UNIX workstation system with hot-standby redundancy. The ES on each workstation acts as an on-line system administrator to diagnose, report, correct, and prevent certain types of hardware and software failures. If a primary station is approaching failure, the ES coordinates the switch-over to a hot-standby secondary workstation. The goal is to discover and solve certain fatal problems early enough to prevent complete system failure from occurring and therefore to enhance system reliability and availability. Test results show that the ES can diagnose all targeted faulty scenarios and take desired actions in a consistent manner regardless of the sequence of the faults. The ES can perform designated system administration tasks about ten times faster than an experienced human operator. Compared with a single workstation system, our hot-standby redundancy system downtime is predicted to be reduced by more than 50 percent by using the ES to command and control the system.

INTRODUCTION

Product reliability and availability are the two most important qualities that ARINC has been pursuing. Currently, a number of systems and products are developed on computer systems running under the UNIX Operating System (OS). Because most off-the-shelf UNIX workstations are general-purpose machines, no fault tolerant mechanisms are built into either the operating system or the hardware. Therefore, hardware or software failures are not automatically detected by the OS utilities. In other words, most UNIX systems do not have any built-in *intelligent* fault diagnostics, fault-correction, or failure-prevention capabilities.

This research project develops a monitor, command, and control ES acting as an on-line system administrator to detect, report, correct, and prevent certain types of hardware and software failures on a UNIX workstation system with hot-standby redundancy. The goal is to use the ES to discover and solve certain fatal problems early enough to prevent system failure from occurring and therefore enhance the system reliability and availability.

The reason we are investigating the potential of using the ES to monitor, command, and control the UNIX workstation is because we need software to perform some of the system administrator's jobs in a real time, automatic fashion. Usually system administrators use *rules of thumb* logic gained through experience to solve problems. ES are the most effective approach to capture and use *rules of thumb* logic to solve problems.

SYSTEM OVERVIEW

As shown in Figure 1, a redundant workstation system consists of two processors—A and B—both powered on, where processor A is the original primary machine and processor B is the secondary. The two processors communicate via StarLan.

One ES resides on each machine, and each ES accomplishes two tasks: (1) process self-checking and hardware management and (2) hot-standby switch-over.

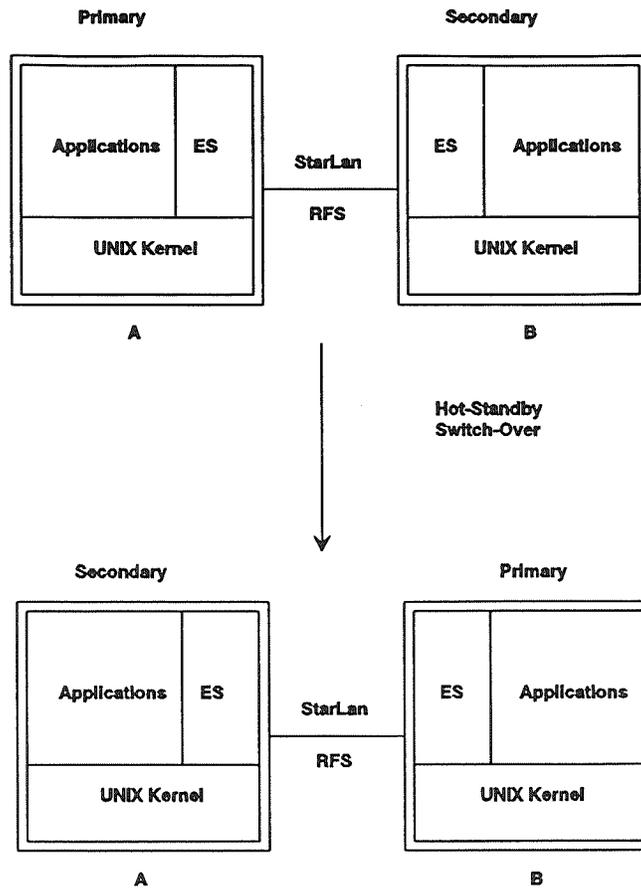


Figure 1. Command and Control Expert System Overview

Processor Self-Checking and Hardware Resource Management

Each ES monitors the health (such as disk space usage, CPU usage, and application process health) of the machine on which it resides. If hardware resource utilization reaches a certain threshold, the ES will first try to correct the situation by taking some actions that a system administrator would take under the same circumstances. If the ES detects problems that it cannot correct or problems that may cause system failure, it will give warnings to the operator. Human intervention is needed in this situation.

Hot-Standby Switch-Over

The two ESs exchange the health information of their host processors with each other. They also monitor the link between the two processors. If the ES on the primary processor detects potential fatal problems, it will send a signal to the secondary processor, and the secondary processor will take over the primary role. Because the switch occurs before the original primary processor failure, there is no system downtime.

The ES uses human heuristic rules to monitor, command, and control the hot-standby redundant system in real time. Unlike other commonly used fault tolerant systems, which take action after the failure occurs, this ES tries to predict failure and take action before failures occur.

HARDWARE AND SOFTWARE SELECTION

Why Use Expert Systems?

A processor usually has multiple processes running simultaneously, has many peripherals, and handles complicated communications. To ensure that the processor operates properly, and that tasks are performed successfully, the processor health status have to be monitored and adjusted continuously. Ideally, problems leading toward processor failure should be discovered and reported, if not solved, before the failure occurs. However, most processors cannot perform self-checking and error-correcting by themselves, and it is unrealistic to have human experts monitor and adjust the processors all the time. Therefore, in most applications, serious problems or even failures occur before human experts issue correction commands.

ES are suitable tools for performing processor monitoring, command, and control for the following reasons:

- Human experts use *rules of thumb* logic gained through experience to solve processor problems; ESs are the most effective tool (in comparison with other types of software) to capture the rules of thumb and use them to solve problems.
- The two research areas of this project, i.e., processor self-tuning and the hot-standby switching, do not necessarily need physical maneuvering. Control commands can be executed by software.
- The processor problems are rather complicated and therefore need specific knowledge and in-depth reasoning to be solved. Conventional programs are not equipped with the features to effectively capture and retrieve knowledge and reasoning.
- Experts agree on solutions. Domain knowledge can be translated into rules and relationships.
- Successful expert systems have been developed in similar applications [11].

Hardware Platform

The target processor for this project is the AT&T/6386, SCSI based workstation, running under UNIX SYSTEM V OS. Two workstations are connected through a StarLan network.

We chose this hardware platform because of its similarity to the ARINC Tower Data Link Service (TDLS). This would allow the results from this research effort to be easily used by TDLS or similar projects. Because UNIX is a machine-independent OS, the research results could also be applied to other UNIX platforms (such as SUN, HP) with little modification.

Software Development Tools

The ES shell chosen for this project is the Cxpert by Software Plus Ltd. Cxpert supports commonly used knowledge representation methods such as *Attributes*, *Frames*, *Procedures* and commonly used inference chaining methods such as *forward chaining* and *backward chaining*. Cxpert also provides query and window display facilities. In addition, Cxpert knowledge representation language (KRL) is fully compatible with C language.

UNIX PROCESSOR SELF-TUNING ES

Requirements Analysis

When application processes run on a UNIX processor, the hardware resources are used by these processes. The health of the processor affects the performance of the applications. For instance, when the CPU is overloaded, the response time of the application process is slower.

This ES prototype shall monitor and maintain the health of the processor in real time in three important areas: disk space, CPU load, and process management. It will solve and report problems in a manner similar to a system administrator, except that it can perform its task 24 hours a day continuously, while humans cannot (because it would be too expensive to be practical).

Knowledge Formalization

The self-tuning (or resource management) ES shall handle the situations stated in the *CONDITION* side of the following IF-THEN statements, and the action the ES takes is on the *ACTION* side of the statements. The following IF-THEN statements are called production rules. They are the body of the ES knowledge base (KB), and they are structured from human heuristics used to handle the same situations.

For disk space management, disk utilization in terms of disk usage percentage is checked every specified time period, as follows:

IF disk usage reaches a certain user-defined high water mark, and

IF file archiving is necessary,

THEN the ES will archive the specified files to tape and then remove them from the disk to create free disk space.

IF the files were archived earlier,

THEN only cleanup is performed.

IF disk usage increases at an abnormally fast rate,

THEN the ES will give the operator a warning.

IF disk usage reaches critical threshold (which means that after this threshold, any writing to the disk has a high probability of failure),

THEN the ES will announce that the processor is approaching a failure condition, and when switch-over is an option, it will initiate switch-over from the primary machine to the secondary machine.

The CPU load is monitored by the ES every specified time period, as follows:

IF the CPU is overloaded,

THEN, to improve system response time for critical processes, certain low priority processes will be terminated.

IF the CPU is still overloaded,

THEN the operator will get a warning.

In many airline or communications related applications, certain application processes must run continuously. The process management part of this ES runs checks every specified period of time:

IF the critical processors are not running,

THEN the critical processes are restarted by the ES automatically.

These are several of the simple but important scenarios that often occur on an applications processor. We studied the logic and procedures a system administrator would take under these scenarios and reformulated human knowledge into an *IF-THEN* format to facilitate ES implementation.

The ES can be expanded to handle more complicated and diversified problems, as long as the human heuristic for those problems can be structured into logical rules.

High-Level Designing

The high-level design for the resource management ES is illustrated in Figure 2. Each block depicts a logical function in the ES, and the arrow indicates data and its flow directions. This block diagram is applicable to each one of the three areas above.

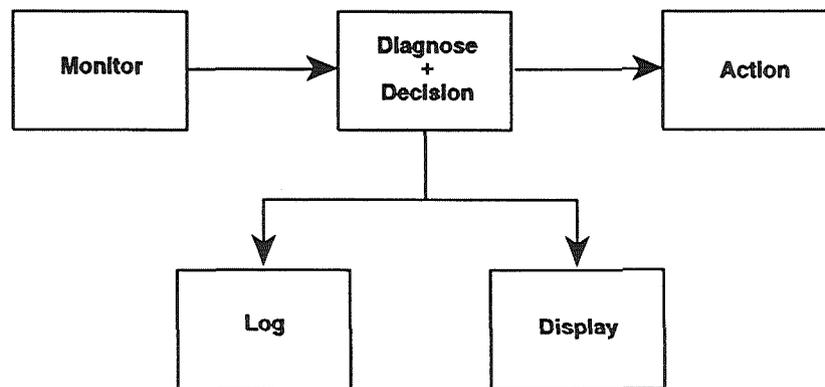


Figure 2. Self-Tuning ES Block Diagram

The *Monitor* (for each one of the areas above) uses UNIX system administration commands to collect resource utilization data. For example, the disk monitor function uses the *df* command to determine free disk space and other information about disk usage. *Monitor* then processes this data so that it can be used by the following processes depicted in the diagram:

- The *Diagnose and Decision* block is essentially a rule-based knowledge base that is formalized from human heuristic to diagnose the health of the resource utilization. This knowledge base fires different rules according to the information from the *Monitor*. The result of a fired rule is an action item, which selects and executes the proper action(s). This block also sends the action item into the log and display blocks.
- The *Action* block is a set of functions that will maintain the health of the processor resource usage.
- The *Log* saves the action item, time stamp, and brief health message onto disk storage.
- The *Display* displays health information and action items on a window system.

Implementation Details

For all three managed areas, the implementation is similar. *Monitor* and *Action* are implemented by C language, *Diagnose and Decision* and *Log* are implemented by Cxpert KRL, and *Display* is implemented by the Cxpert Hyperwindow system.

Performance Test

Test-run results show that the hardware resource management ES can handle all the fault scenarios described above and execute corrective actions regardless of the sequence of faults. The ES behavior is fully predictable. The one-round (i.e., a visit to the three targeted areas once) ES execution time is less than 15 seconds, which is much faster than comparable human action. A more sophisticated ES may have a longer one-round execution time, but it will still be much faster than manual operation. The performance speed can be improved by using C or C++ language for coding. In addition, operator error can be avoided because all the proper actions that the ES takes are pre-coded into the ES and tested.

HOT-STANDBY SWITCHING ES

Requirements Analysis

As shown in Figure 2, the hot-standby redundant system consists of two identical UNIX processors linked by StarLan Ethernet for peer-to-peer communications. When the system is first started up, one processor is designated to be primary and the other to be secondary. Hot-standby implies both machines are powered, and both are accepting and processing the same input. However, only the results of the primary processor are used as system output. Hot-standby switching requires that when the primary processor fails, the primary role is switched to the secondary processor within a specified period of time. This period of time must be short enough so that the *outside world* will not be affected by the processor failure.

Many mechanisms can be used to initiate hot-standby switching[7]. The ES uses the following mechanism:

Each processor uses the self-tuning ES to check its own health and sends its health information to its peer processor via the StarLan link. If the primary processor detects its own failure, it stops its normal operation and sends a signal to the secondary processor. The secondary configures itself to be the primary processor and picks up the operation where the other machine stopped.

Knowledge Formalization

The production rules for the hot-standby switching ESs in the two machines are slightly different because of the different roles they play. The two processors use the self-tuning ES to check self-health and send this information to their peer processor via the link. In the meantime, each processor tries to read the health of the other machine.

In the primary machine-hosted ES, the rules are as follows:

IF the StarLan link is healthy, and

 IF the secondary processor is healthy, and

 IF the primary is going to fail (in the prototype, only disk problems are considered as fatal),

 THEN the primary stops its output, announces its failure to the secondary, and sends out a switch-over command.

 IF the secondary processor is NOT healthy,

 THEN the primary processor will report that its peer is dead; no switch-over is allowed under this circumstance.

IF the link is dead,

THEN the primary machine will report this; no switch-over is allowed under this circumstance, either.

In the secondary processor, the rules are as follows:

IF the link is healthy, and

IF the secondary is healthy, and

IF the primary is dead, and

IF the primary sends a signal for switch-over,
THEN the secondary processor reconfigures itself to be primary and announces that it is primary now.

IF there is no switch signal from the primary,
THEN the secondary will announce this fact, but no switch-over occurs. (Note: this rule is intended to avoid a race condition between the two machines for the primary role and confuse or upset the *outside world*. A more elaborate rule set will allow the secondary to take over without the primary switch signal.)

IF the secondary is dead,
THEN it will notify the primary, and no switch signal will be sent to the secondary in case of primary failure.

IF the link is dead,
THEN no switch-over is allowed.

Listed above are some simple rules to control and coordinate the hot-standby switch-over process. More elaborate rules are necessary for a real operation context. These rules are for prototyping purposes only. The ES can be expanded to handle more complicated situations.

High-Level Design

Figure 3 illustrates the functional blocks in the hot-standby switching ES. The functionality of each block is explained by the block name. The ES high-level design is the same for both ESs, except that the primary machine has the *Send Switch* while the secondary machine has the *Receive Switch*.

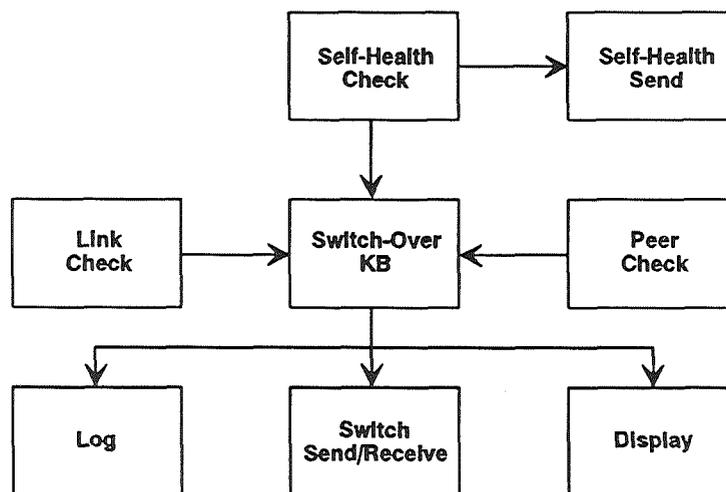


Figure 3. Hot-Standby Switching ES Block Diagram

Prototype Console Display

Figure 4 shows the primary and secondary machine console display. In the *DISK* window, machine disk space usage is given as a percentage. The ES action for self-tuning is also shown in the window. In the *CPU* window, average run queue size, average queue occupancy data, and the ES action for process control (such as start/kill process) are displayed. The *PEER* window shows whether the peer processor is healthy. The *LINK* window shows whether the RFS is still running on StarLan. The *SWITCH* window shows whether the host machine is primary or secondary. The *LOG* window logs the time, problem, and action taken by the ES.

There are also colors associated with the windows to indicate the status of the related area. If a critical situation occurs in a certain area, the corresponding window will turn red; if a warning occurs, the window will turn white, etc.

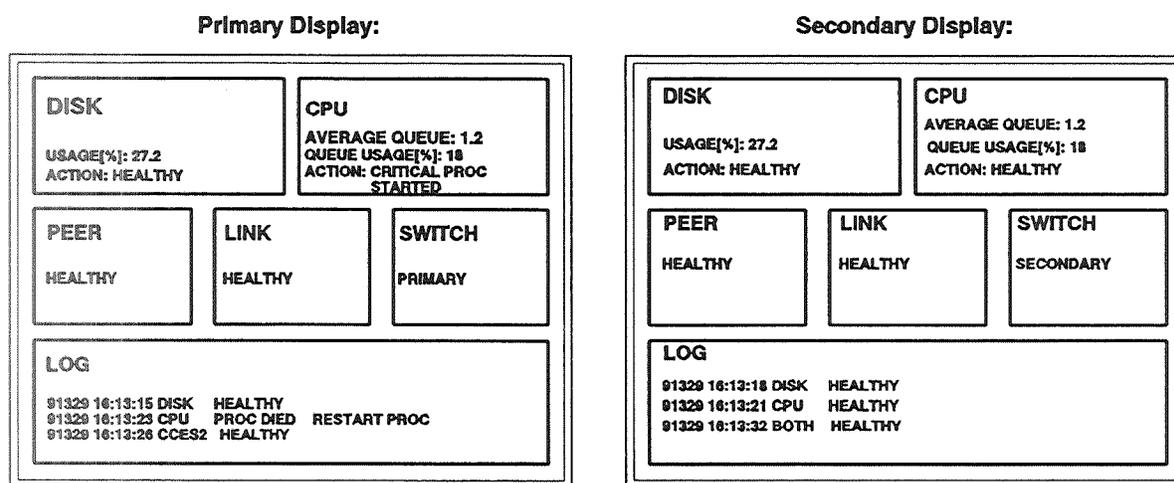


Figure 4. Prototype Console Display

Implementation Details

The *Switch-Over Knowledge Base (KB)* is implemented by Expert KRL, and the rest of the blocks are achieved by C functions.

The hot-standby switching ES incorporated the ES in the section titled "UNIX Processor Self-Tuning ES" for processor self-checking and self-tuning.

In this ES, in addition to the knowledge base, the technical challenge is the Inter-Process Communication (IPC). Regular IPC facilities such as message queues, unnamed pipes, and semaphores cannot communicate across the network, and the named pipes cannot satisfy the independency requirement between the two processors. This ES implementation uses a unique and innovative method that combines the StarLan Remote File Sharing (RFS) facility with file and record-locking techniques to accomplish the peer-to-peer communication across the network. This implementation will allow the two processors to operate independently if the link or one of the processors is dead. It will also allow the two processors to operate cooperatively if the link is healthy and both of the processors are running.

Performance Test

In addition to all the faulty scenarios presented in the section titled "UNIX Processor Self-Tuning ES", tests are generated to simulate link failure, primary failure, secondary failure, and both processor failures. The ES

performs correctly under these circumstances. The one-round execution time (which is the time that the ES takes to check the three self-tuning areas, the peer, and the link once) is less than 30 seconds, which means that the secondary will take over the primary function in less than 30 seconds in case of primary failure. A more sophisticated ES may have a longer one-round time, but performance speed can be improved by using C or C++ language to avoid Cxpert KRL overhead costs.

CONCLUSIONS

Command and Control Expert System (CCES) Potential Applications

The commercial application of the ARINC CC ES is very wide, because its knowledge base captures general UNIX system administration knowledge, and the rules that CCES uses are applicable to any UNIX system with little or no design change. Only the detailed implementation may vary for different platforms. For example, in the air traffic control industry, this ES has potential applications in a number of FAA directed services, the TDLS is one example.

Conclusions

The UNIX processor self-tuning ES prototype allows the processor to have better hardware resource management and, therefore, better performance and less chance of failure.

The hot-standby switching ES prototype provides coordination, command, and control to the switch-over process and, therefore, reduces system downtime and improves system reliability and availability.

REFERENCES

- [1] A Guide to Expert Systems, Waterman, Addison-Wesley, 1986
- [2] "Applying Systems Analysis Techniques to Knowledge Engineering", Expert Systems, Vol. 7, No. 2, May 1990, Swaffield, G. and Knight, B.
- [3] "Architecture of Fault-Tolerant Computers", IEEE Computer, Siewiorek, D. P., August 1984
- [4] "Artificial Intelligence Technologies for Real-Time and Object-Oriented Applications," Electrical Communication, Vol. 62, No. 3/4, Barachini, F., 1988
- [5] AT&T StarGROUP Software Reference Set, AT&T, 1988
- [6] Building Expert Systems, Hayes-Roth and Waterman, Addison-Wesley, 1983
- [7] Design and Analysis of Fault Tolerant Digital Systems, Johnson, J. W., Addison-Wesley, 1989
- [8] "Design for Ultrahigh Availability: The Unix RTR Operating System," IEEE Computer, Wallace, J. J. and Barnes, W. W., AT&T Bell Laboratories, August 1984
- [9] "Expert Systems Making Quiet Inroads into Networking Applications," Networking Management, May 1991
- [10] "The Real-Time Expert," BYTE, Laffet, T.J., January 1991
- [11] "YES/MVS: A Continuous Real Time Expert System," Proceedings AAAI-84, Griesmer, J. H., <et al>; IBM, 1984
- [12] "New Controls for Air Traffic," IEEE Spectrum, February 1991

528-61
150498
P-11

The Generic Spacecraft Analyst Assistant (GenSAA): A Tool for Developing Graphical Expert Systems

N 9 3 - 2 5 5 8 9

Peter M. Hughes
Computer Engineer/GenSAA Project Manager
Software and Automation Systems Branch (Code 522)
NASA/Goddard Space Flight Center
Greenbelt, Maryland 20771

Abstract

During numerous contacts with a satellite each day, spacecraft analysts must closely monitor real time data watching for combinations of telemetry parameter values, trends, and other indications that may signify a problem or failure. As the satellites become more complex and the number of data items increases, this task is becoming increasingly difficult for humans to perform at acceptable performance levels. At the NASA Goddard Space Flight Center, fault-isolation expert systems are in operation supporting this data monitoring task. Based on the lessons learned during these initial efforts in expert system automation, a new domain-specific expert system development tool named the Generic Spacecraft Analyst Assistant (GenSAA), is being developed to facilitate the rapid development and reuse of real-time expert systems to serve as fault-isolation assistants for spacecraft analysts. Although initially domain-specific in nature, this powerful tool will readily support the development of highly graphical expert systems for data monitoring purposes throughout the space and commercial industry.

Introduction

NASA's Earth-orbiting scientific satellites are becoming increasingly sophisticated. They are operated by highly trained personnel in the mission's Payload Operations Control Center (POCC). Currently at the Goddard Space Flight Center (GSFC), missions utilize either a dedicated control center (e.g., LANDSAT and the Hubble Space Telescope) or share resources in the Multi-Satellite Operations Control Center (e.g., Cosmic Background Explorer and the Gamma Ray Observatory). In either case, POCC personnel called Flight Operations Analysts (FOAs), are responsible for the proper command, control, health, and safety of the satellite¹ [3].

The satellite control centers operate round-the-clock throughout the lifetime of the spacecraft. There are typically multiple real-time communications events daily with each satellite. During these events, the FOAs must:

- establish and maintain the telecommunications link with the spacecraft,
- monitor the spacecraft's health and safety,
- send commands or command loads to the satellite for on-board execution,
- oversee the transfer of the scientific data from the on-board tape recorders to ground systems for processing and analysis, and
- manage spacecraft resources (including on-board memory, batteries, and tape recorders).

To accomplish these activities, the analyst must thoroughly understand the operation of the spacecraft and ground systems and continuously monitor the current state of operations as indicated by telemetry parameters displayed on the POCC consoles. During an event, the analyst typically monitors hundreds of telemetry parameter values on multiple display pages that may be updated several times per second. Such large volumes of low-level information can overwhelm analysts and disrupt their ability to identify and resolve conflicting constraints. They may soon be unable to consistently monitor all of the information available. The need to automate some data monitoring functions is apparent.

Expert system technology is proving to be effective in automating some spacecraft monitoring functions. This paper first summarizes CLEAR, the first real-time spacecraft monitoring expert system deployed at GSFC. The paper then reviews several lessons learned from CLEAR and other monitoring and fault isolation expert system projects undertaken at GSFC, thereby establishing the foundation of a domain-specific expert system development tool called the Generic Spacecraft Analyst Assistant (GenSAA). This new tool will be introduced followed by a discussion of its capabilities, architecture and benefits, its potential uses in industry, and the approach for adapting this tool to other domains.

1. This paper is based on a previous publication by Hughes, P. & Luczak, E. (October, 1991), reference 3.

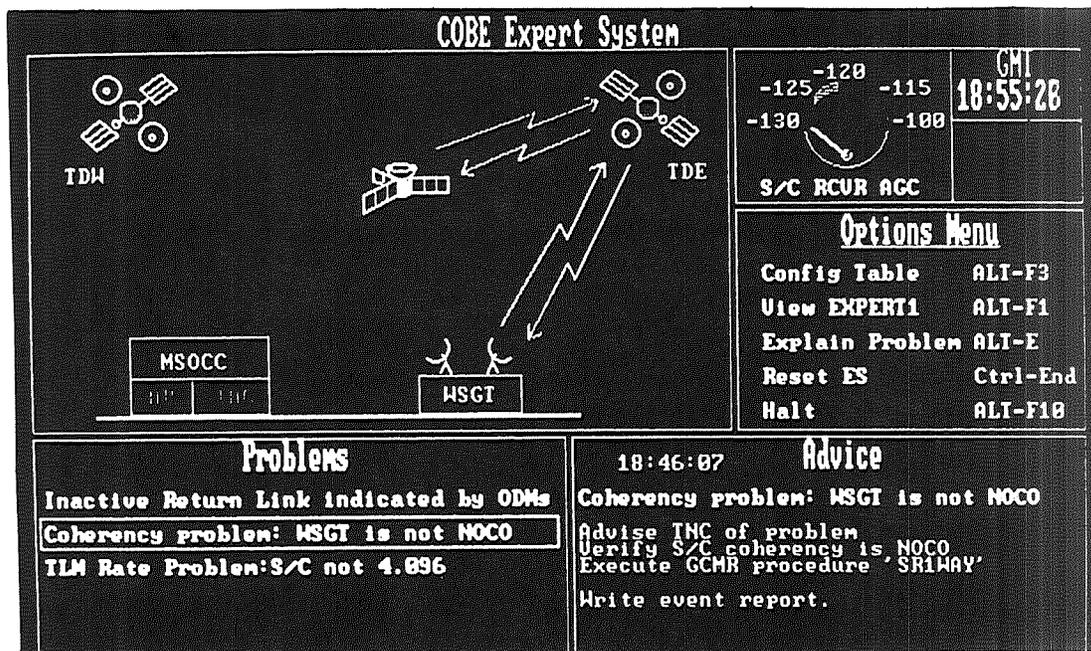


Figure 1. Black & White Photo of the CLEAR User Interface

CLEAR: Breaking New Ground

The Communications Link Expert Assistance Resource (CLEAR) is the first operational expert system at GSFC that automates one of the spacecraft analyst's tasks [2]. It is a fault-isolation expert system that supports real-time operations in the POCC for the Cosmic Background Explorer (COBE) mission. CLEAR monitors the communications link between COBE and the Tracking and Data Relay Satellite (TDRS), alerts the analyst to any problems, and offers advice on how to correct them.

CLEAR is a forward chaining, rule-based system that operates in the COBE POCC. It monitors over 100 real-time performance parameters that represent the condition and operation of the spacecraft's communications with the relay satellite. Using this information, together with knowledge of TDRS operations, COBE's on-board communications system and the expected configuration of the scheduled event, CLEAR accurately portrays the status of the communications link.

Textual and graphical information about the condition of the COBE/TDRS/ground communications links is displayed in a tiled-window format (Figure 1). A graphics window displays the elements of the communications network from the COBE Spacecraft to the POCC; green lines represent healthy links between elements. When the performance parameters indicate that a communications link or processing system is degrading or down, the associated line or icon turns yellow or red, respectively. The display enables analysts to assess the current status of the communications event in a quick glance.

When CLEAR isolates a problem, a short description of the problem is displayed in the "problems" window. If multiple problems are found, the problem descriptions are ranked and displayed in descending order of criticality. CLEAR suggests analyst actions to correct the problem; however, the system does not take any corrective action itself.

To further assist the analyst and to provide support for its advice, the CLEAR system provides an explanation facility. When the analyst selects a problem displayed in the problems window, CLEAR provides a detailed explanation of why the expert system believes that the problem exists.

CLEAR has approximately 165 rules and isolates approximately 75 different problems. The types of problems include: non-reception of data within the control center (system or communication problems, or data reporting not activated); misconfigurations between the COBE POCC and the TDRS ground station (coherency/non-coherency, doppler compensation on/off, power mode, actual TDRS in use, antennae configurations); discrepancies in telemetry rate or format; inactive or non-locked links; and degrading or critical signal strength situations [6].

CLEAR operates on any of the seven PC/AT-class workstations that are used for console operations in the POCC. It is written in the 'C' language and uses the 'C' Language Integrated Production System (CLIPS) and a custom-developed graphics library.

The CLEAR Expert System has supported the COBE Flight Operations Team since launch in November 1989. It is used to monitor nearly all of the TDRS supports (COBE occasionally communicates directly to the Wallops ground station) and is regarded as the fault-isolation "expert" for the COBE/TDRS telecommunications link. CLEAR represents a successful attempt to automate a control center function using an expert system. It has been adapted for the Gamma Ray Observatory and was utilized during early orbit.

Lessons Learned

Several important lessons have been learned from the experience gained in developing and operating CLEAR [2]. Key lessons have also been learned from other monitoring and fault isolation expert systems developed recently at GSFC, including the Ranging Equipment Diagnostic Expert System (REDEX) [5], and other systems. These lessons learned have strongly influenced the definition of GenSAA.

- Production rules effectively represent analysts' knowledge for automating fault-isolation in spacecraft operation. The rule-based method of knowledge representation has proven to be quite powerful for fault-isolation in spacecraft operations. Production rules provide a direct method of encoding the fault-isolation knowledge of spacecraft analysts; the if-then structure closely parallels the stimulus-response behavior that they develop through extensive training. This knowledge can be translated smoothly into rule form. The development of CLEAR would have taken much longer using conventional, non-rule-based programming techniques.

- Knowledge engineering is an iterative, time-consuming process. Early in CLEAR's development, the primary concern was the perceived difficulty of the knowledge acquisition effort. However, the knowledge engineering task was found to be relatively straightforward, albeit time-consuming. The development of the rule base was a lengthy process due to the interactive nature of the knowledge acquisition. Basically, the expert would describe a specific piece of knowledge to the "knowledge engineer" who would attempt to transcribe it into a rule and pass it back to the expert for validation. When the rule accurately represented the piece of knowledge (which usually took multiple iterations between the expert and the knowledge engineer), it was passed to the test team for formal testing, and then, finally, released for operational use.

The involvement of various players in this process resulted in long turnaround times from the point at which a piece of knowledge was determined to be important until it was translated into a rule and placed into operation. It was recognized that an integrated tool that simplified modifications to the knowledge base or user interface could significantly accelerate this process.

- Allow the domain expert to participate in the rule formation. The CLEAR development team learned that the knowledge structure of the fault-isolation process employed by the FOAs is shallow (i.e., the identification of a problem generally does not rely on the identification of other subproblems, and so on). Most of the problems identified by the analysts were discrete problems that seldom overlapped other problems. Conflicts between rules were minimal; this simplified testing, verification, and validation of the rulebase.

The participation of the analyst in knowledge acquisition and translation has many advantages. First, it can reduce the knowledge translation time and, more importantly, reduce knowledge translation errors that occur when a knowledge engineer formulates rules based on the knowledge extracted from documentation or interviews with the domain expert. Second, the verification and validation of the knowledge will be facilitated since the expert will better comprehend the rulebase. Third, the in-depth understanding of the knowledge base and its capabilities is likely to result in a higher degree of user confidence in the system thereby ensuring user acceptance.

- Expect to fine-tune the expert system after it becomes operational. For CLEAR, the rule-based method of knowledge representation has provided the flexibility to easily adapt the knowledge base to unforeseen changes in the operational behavior of the spacecraft. For example, even though the operational nature of COBE was fairly accurately understood by the design engineers and flight operations team before the launch, slight behavioral variations and complications arose once the spacecraft was in orbit. Although the FOAs were able to adjust to such variations quickly, some of the ground systems required complex software modifications. However, the changes required to CLEAR's rule-base were simple and quickly implemented.

After modification, CLEAR provided consistent operational assistance. It is important to provide the capability to modify an operational expert system in a controlled, yet straightforward way.

- Don't underestimate the integration process. One of the most valuable lessons learned is that while prototypes can often be developed rapidly, operational expert systems require considerable effort. A major factor in this effort is the difficulty of interfacing the expert system to the data source.

The CLEAR development team learned that most of the development time for the operational system was spent on issues not directly related to the construction of the knowledge base. A surprising amount of effort focused on the integration of the expert system with the data source and graphics display system. This required in-depth programming knowledge of the interfacing systems and the ability to troubleshoot problems within them. Provide tools to simplify the complicated task of integrating the expert system with the interfacing systems and, if possible, reuse any interface code developed for a similar (expert) system.

- Don't neglect the user-interface. The human-computer interface is frequently the most underdeveloped component of an expert system. An effective user interface is inviting, comprehensible, credible, and simple to operate. To make it inviting, simplify the display layout and present only that information needed to efficiently perform the task. Graphics can greatly enhance the visual communications of a system; capitalize on their expressive power to provide system output that can be assimilated quickly and accurately.

- Use graphical diagrams to illustrate the system being monitored. Users have responded very positively to the use of schematic displays that graphically represent system status and fault locations. Analysts and technicians usually learn about the systems they monitor by studying system block diagrams in training classes and reference manuals. By using similar block diagram displays, a monitoring expert system can present status to the user in a familiar and intuitive format. Color coding of status conditions on such displays has been found to be an effective way to present succinct status summaries.

- Make the system easy to operate. Operation of the expert system should be simple enough that the user can concentrate on the problem, not on how to operate the system. The following techniques were applied in CLEAR and REDEX to simplify operation:

- Reduce user input to a minimum. CLEAR operates in a highly autonomous mode; no user input is required after system initialization. CLEAR has been well-accepted by its users, partially because it operates as an independent intelligent assistant, allowing the spacecraft analyst to focus on other responsibilities during real-time satellite contacts.

- Use hypergraphic techniques. These techniques [1] enable the expert system user to quickly select and display desired diagrams by clicking on link buttons that appear on each diagram. Links can be used to create diagram hierarchies, off-page connections, diagram sequences, and other structures.

These lessons learned have actuated the definition and development of GenSAA.

GenSAA

Overview

GenSAA is an advanced tool that will enable spacecraft analysts to rapidly build simple yet highly graphical expert systems that are capable of performing real-time spacecraft monitoring and fault isolation functions. Expert systems built using GenSAA will assist spacecraft analysts during real-time operations in spacecraft control centers. The use of GenSAA will reduce the development time and cost for new expert systems in this domain. GenSAA will allow graphical displays and fault-isolation knowledge to be reused from mission to mission.

Expert systems developed with GenSAA will have the following characteristics:

- Easily created and modified— The process for developing specific expert systems using GenSAA will be straightforward enough that it can be performed by trained spacecraft analysts on the flight operations team. No compilation step is necessary before executing the expert system.

- Rule-based— GenSAA will support the use of rules to represent spacecraft and payload monitoring and fault isolation knowledge. Rule-based representations are easily learned and can be used to describe many of the reasoning processes used by spacecraft analysts.

- Highly graphical— The GenSAA operational user interface will support both textual and graphical presentations of health and status information and fault isolation conclusions. GenSAA user interfaces are built with the GenSAA WorkBench which uses the X-window toolkit and the Motif widget set. Hyperlink

techniques will be supported to simplify navigation between GenSAA windows.

- Transparently interfaced with the data source— Initially, GenSAA will be used to create expert systems that will support analysts in spacecraft control centers that use the new Transportable Payload Operations Control Center (TPOCC) architecture. TPOCC is a new Unix-based control center system architecture that will be used on many new spacecraft missions at GSFC. GenSAA will be adaptable to also support non-TPOCC data interfaces.

- Real time— GenSAA expert systems will be driven by real time spacecraft telemetry that indicate the current status of the spacecraft and its operation.

GenSAA is being developed as a generic tool to support the development of expert systems in any TPOCC-based control center such as SAMPEX, Wind/Polar, SWAS, SOHO, and others. However the initial use of GenSAA will be targeted for the SMEX and ISTEP series of missions. SAMPEX flight operations team members have expressed a need for a tool like GenSAA, and the launch timeframe for SAMPEX, the first SMEX mission, is compatible with the GenSAA development schedule.

GenSAA Architecture

GenSAA is an advanced, domain-specific tool for developing spacecraft control center expert systems. It is analogous to many commercial expert system shells because it facilitates the development of specific expert systems. However, GenSAA is tailored to the specific requirements of spacecraft analyst assistant expert systems in TPOCC control centers.

GenSAA operates in the TPOCC environment and shares many of TPOCC's architectural features. The TPOCC architecture is based on distributed processing, industry standards, and commercial hardware and software components. It employs the client/server model of distributed processing, the Network File System (NFS) protocol for transparent network access to files, and the X Window System (X.11) with the Motif library and window manager for the graphical operator interface. A TPOCC configuration consists of a small set of specialized front-end processors and Unix-based workstations on an Ethernet network using the TCP/IP network protocol. GenSAA operates in this environment.

Figure 2 shows the major elements of GenSAA. They are divided into two sets: the GenSAA Workbench and the GenSAA Runtime Environment. The Workbench is used in an off-line mode to create a specific GenSAA Expert System and the Runtime Environment is used to execute the GenSAA Expert System to support real-time operations in a spacecraft control center. These elements are described in the sections below.

The GenSAA Workbench

The GenSAA Workbench is composed of three utilities that enable a spacecraft analyst to create a GenSAA expert system. A GenSAA expert system is a specific expert system that performs real-time monitoring and

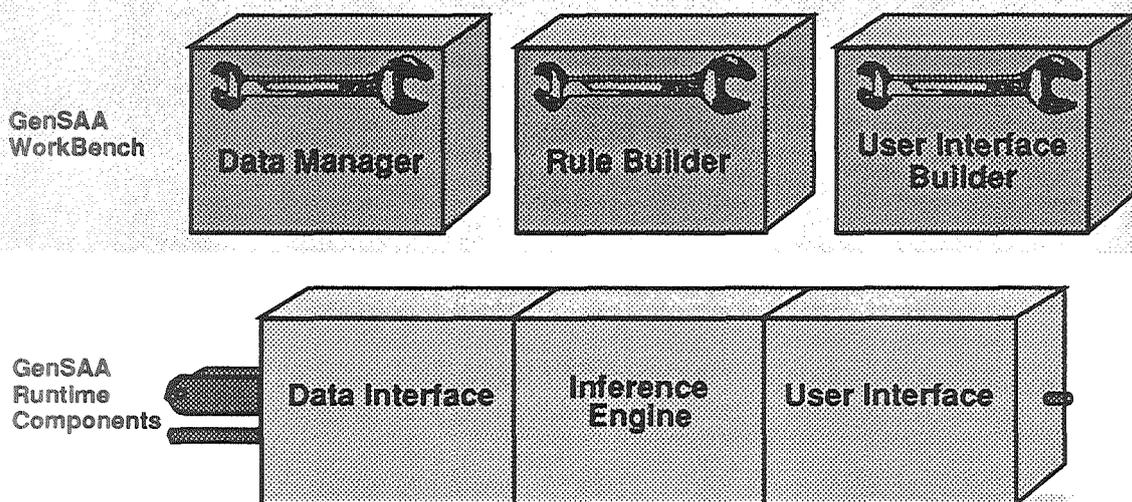


Figure 2. The Elements of GenSAA

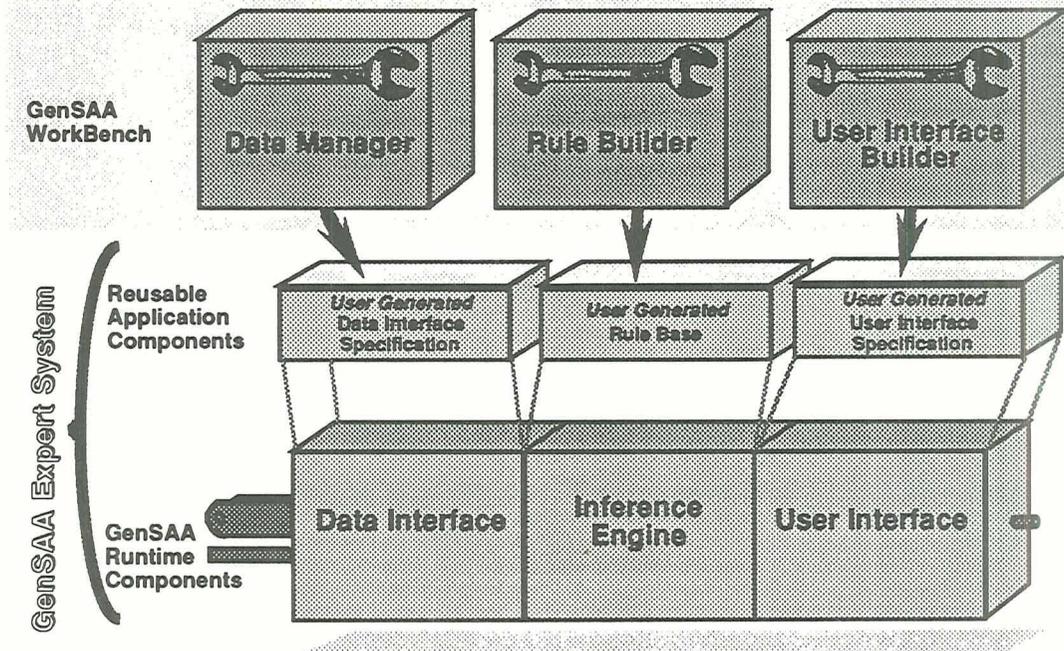


Figure 3. Creating a GenSAA Expert System

fault isolation functions in a TPOCC spacecraft control center.

The GenSAA Workbench will operate in an off-line mode on a Unix workstation. A GenSAA expert system is created by defining the expert system's runtime specifications using the GenSAA Workbench. Figure 3 illustrates that these specifications, called Reusable Application Components, together with the GenSAA Runtime Components, compose a GenSAA Expert System. The GenSAA Workbench utilities are as follows:

- **Data Manager**– This utility is used to create the Data Interface Specification for a GenSAA expert system. The Data Interface Specification defines four types of data that are used by the GenSAA expert system during real-time operations:

- *Mission data*– Mission data variables represent real-time status of the monitored spacecraft and related ground support systems. (mission variables are sometimes called telemetry mnemonics.) Values for these variables are received and updated during spacecraft operation periods from the TPOCC Data Server process, which is part of the TPOCC software. Using the Data Manager, the GenSAA Workbench user selects the mission variables needed for the expert system being created from a list of all the mission variables available from the TPOCC Data Server. Values for only those variables selected will be received by the expert system during run-time.

- *User-defined data*– User-defined data variables represent expected operating modes and equipment configurations. For example, a user-defined data variable might represent the setting of a switch that determines which of two redundant components is to be used. Values for these variables are entered by the spacecraft analyst during spacecraft operation periods.

- *Inferred data*– Inferred data variables represent conclusions inferred by rules in the rule base. For example, an inferred data variable might represent the health or fault status of a component in a spacecraft subsystem. (The name of an inferred data variable together with its current value is called an inferred fact.) Values are assigned to these variables by actions executed in the “then” part of rules that fire.

- *Externally Generated GenSAA (EGG) data*– EGG data consists of Inferred and User-Defined data which is identified by the user as being “public”. These data are routed to the GenSAA Data Server to make them available to any process requesting them by name. For example, one GenSAA expert system may require information about the status of a subsystem which is being monitored by another GenSAA expert system. Such inter-expert system communication is conducted through EGG data.

- **Rule Builder**– This utility is used to create the rule base for a GenSAA expert system. The rule base is a set of expert system rules in “condition-action” (“if - then”) format that may infer new facts based on currently asserted facts. The inference engine controls the firing of rules in the rule base during execution of the GenSAA expert system.

During run-time, if all the conditions of a rule are satisfied, then the rule fires and all its actions are performed. Conditions can be constructed using the mission, user-defined, and inferred data variables specified with the Data Manager. Actions may include: asserting/retracting an inferred fact, enabling/disabling a rule or ruleset, performing a mathematical calculation, and displaying text messages on the user interface. Templates are provided for specifying conditions and actions thereby allowing a user to build rules quickly using drop-and-drag techniques.

- **User Interface Builder**– This utility is used to create the User Interface Specification for a GenSAA expert system. The User Interface Specification defines the user interface windows and the layout and behavior of the graphical objects that comprise the operational user interface of the GenSAA expert system.

The Workbench user can use a variety of X-toolkit and Motif widgets, including pushbuttons, option menus, scrolling text lists, user-created graphical icons, and data-driven objects such as meters and gauges. The designer constructs a user interface by selecting graphical objects from a palette or drawing them with the graphics tools provided and placing them wherever desired. Lines can be drawn using connector items to create animated schematic diagrams. The Workbench user can associate each graphical object with a mission, user-defined, inferred data, or EGG variable, and specify how changes in the value of the variable will affect the presentation of the item. Characteristics of a graphical object's behavior that can change based on the value of its associated data variable include its color, the icon displayed, and the position of the dynamic portion of a data-driven object. Simple drawing editors are provided to allow the creation of new graphical icons. Any graphical object can also be specified to be a hyperlink button; clicking on such a button during run-time can cause a window to be displayed, or cause an informational pop-up window to appear.

The GenSAA Workbench utilities are highly interoperable and use a graphical, direct-manipulation method of interaction (commonly referred as "point-and-select" or "drag-and-drop") to facilitate use. For example, when using the Data Manager, the user may select a given mission mnemonic to be included in the Data Interface Specification. Later, when using the Rule Builder, the user can drag the mnemonic from the Data Manager into a condition of a rule. Similarly, when using the User Interface Builder, the user can drag a GenSAA data variable onto a graphical item in the user interface to associate the variable with the graphical object. This pointing technique prevents keyboard mistakes and is faster than typing.

In Figure 3, the outputs of the GenSAA Workbench utilities are described as reusable application components. These components will be placed in a library so that they can be reused in creating the specifications for new GenSAA expert systems. Operations like cut and paste will be available to allow portions of previously created specifications to be used in constructing a new expert system.

GenSAA Runtime Environment

The elements of the GenSAA Runtime Environment are called the GenSAA Runtime Components; they are used without change in each GenSAA expert system. They control the operation of a GenSAA expert system during its execution in a TPOCC control center. They read the Data Interface Specification, Rule Base, and User Interface Specification files to determine the specific behavior of the GenSAA expert system. The GenSAA Runtime Framework is implemented as a pair of Unix processes that communicate with one another via message queues. Their functions are as follows:

- **User Interface Process**– This component manages the user interface of the GenSAA Expert system. It displays user interface windows that contain both text and graphics. Color is used to enhance the display of state data. Data-driven display objects are associated with telemetry values received from the TPOCC data server and inferred facts and conclusions received from the Inference Engine. In response to user inputs that include hypertext button events, the User Interface displays selected graphics windows, help text, and other informational text. The user interface windows, data-driven objects, and interaction objects are defined in the User Interface Specification that was generated by the GenSAA User Interface Builder .

- **Data Interface Subsystem**– This sub-element of the User Interface Process requests telemetry from the TPOCC Data Server, as specified in the Data Interface Specification. It formats the real-time data it receives and makes it available to the Inference Engine and User Interface components.

- **Inference Engine Process**– This component manages the firing of rules in the rule base. A rule is fired when all its conditions are satisfied; the conditions will often involve the current values of telemetry, user-defined, and inferred data variables. Inferred facts and messages may be sent to the User Interface

component and displayed to the FOA as defined in the User Interface Specification. NASA's 'C' Language Integrated Production System (CLIPS) inference engine forms the core of this component.

The operational interface with the FOA will typically include color schematics and animated data-driven objects (such as rotating meters, sliding bar graphs, and toggle switches) that graphically display the dynamic values of telemetry data, user-defined data, and inferred conditions. The user interface will also typically contain hypertext and hypergraphic links to make it easy for the spacecraft analyst to quickly display graphics windows.

Figure 4 shows a completed GenSAA expert system in operation. A GenSAA expert system will execute on a Unix workstation in the control center. A dedicated Unix workstation is not required, i.e., a GenSAA expert system can execute on the same workstation as other TPOCC processes. However, to avoid potential performance impacts, the initial GenSAA expert system will reside on a dedicated Unix workstation in the SMEX POCC connected to the TPOCC Local Area Network.

During operation, a GenSAA expert system will interface to the TPOCC software via the TPOCC Data Server process. This interface will use the standard external interface conventions defined by the TPOCC Project. For example, a GenSAA expert system simply submits a request to the TPOCC Data Server process for the telemetry items that were specified in the expert system. No additional data or commands are sent to any of the TPOCC processes.

Implementation

The GenSAA development team is utilizing a spiral development approach in which two prototypes and an operational system will be implemented. The first prototype, called the 'Proof-of-Concept Prototype', investigated the basic concepts of GenSAA and was completed in August, 1990. The second, called the 'Functional Prototype', demonstrated the functional characteristics of the GenSAA Workbench. This prototype was completed in October, 1991 and was used to assess and refine the functional requirements for the operational system. The operational version is scheduled for release in mid 1993.

GenSAA will be implemented in C++ using an object-oriented design. This approach has been selected because of the following four benefits: First, it will allow the reuse of an existing class library developed at GSFC (Code 522) for the rapid development of software. Second, it will promote modularity and ease integration of the software components that will comprise GenSAA. Third, it will allow the core modules of

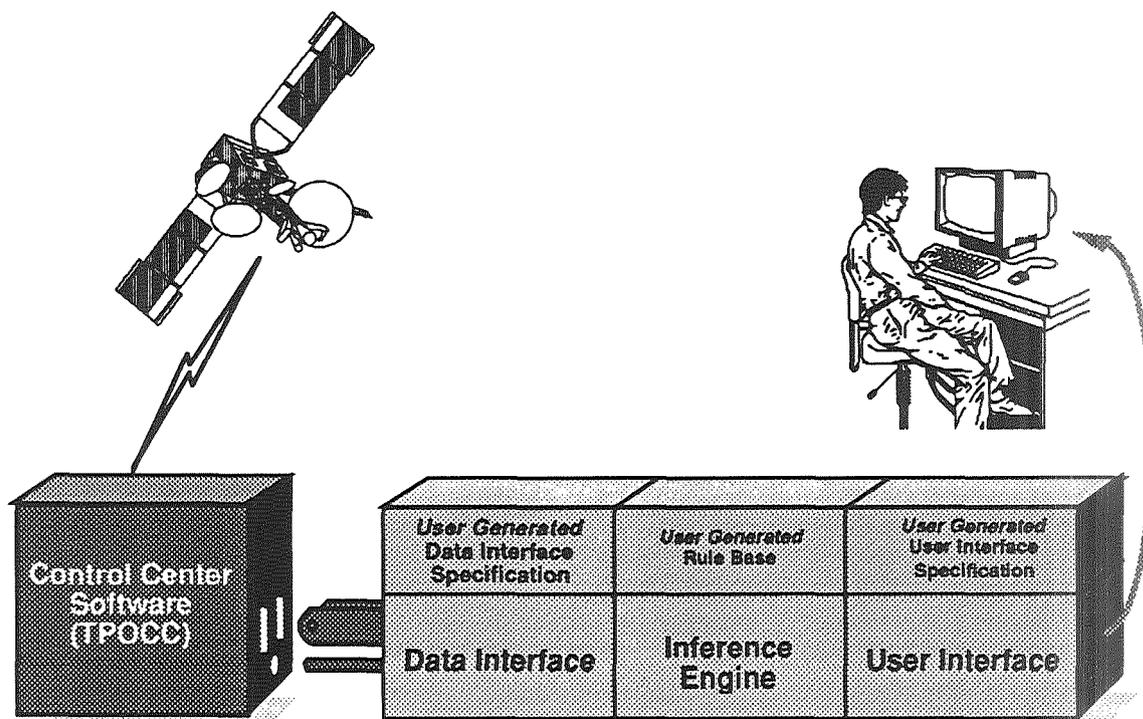


Figure 4. A GenSAA Expert System In Operation

GenSAA to be implemented so that the system can be extended for future missions or industrial use without major design changes or extensive recoding. Fourth, the GenSAA development team is optimistic that the object-oriented approach will facilitate maintenance of this system.

Multiple GenSAA Expert Systems

GenSAA expert systems are intended to be relatively simple expert systems with small rule bases that are typically developed by a single analyst. A typical GenSAA expert system might monitor and isolate faults for one subsystem on board a spacecraft. To handle more complex monitoring situations, involving, for example, several spacecraft subsystems, multiple GenSAA expert systems can be built each responsible for a discrete subsystem or function. During operation, these expert systems would execute concurrently and could share key conclusions with one another using a "publish-and-subscribe" model of communicating.

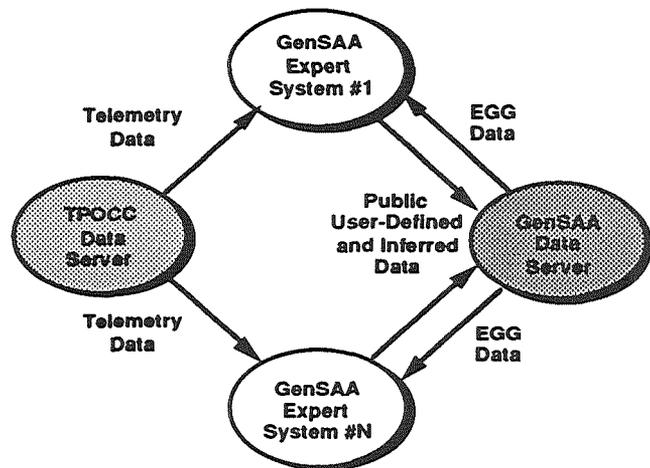


Figure 5. Sharing data among multiple GenSAA Expert Systems

To perform the publish-and-subscribe method of information sharing, a fourth component of the GenSAA Runtime Environment, the GenSAA Data Server, is used to serve as a central repository to which GenSAA Expert Systems can "publish" information and from which other "subscribing" GenSAA Expert Systems can receive the information when published. As shown in Figure 5, the GenSAA Data Server is a Unix process that can receive a real-time stream of user-defined and inferred data variable updates from any GenSAA expert system. The GenSAA Data Server distributes the data to any GenSAA expert system that has requested it. A given GenSAA expert system only receives those variables it specifically requested (subscribed). The data received by a GenSAA expert system from the GenSAA Data Server is called externally generated GenSAA (EGG) data. A GenSAA expert system receives EGG data via its Data Interface component in exactly the same way as it receives telemetry data from the TPOCC Data Server.

Within a GenSAA expert system, EGG data can be used in the conditions of rules, and can be associated with display items in exactly the same way as mission, user-defined, and inferred data. The Workbench supports the selection of EGG data as a fourth variable type. The Workbench also allows any local user-defined or inferred data to be specified as public, to cause it to be sent to the GenSAA Data Server, and thereby shared with other GenSAA expert systems.

Benefits of GenSAA

The following benefits are expected to be realized by using GenSAA to build spacecraft monitoring expert systems for future NASA missions:

- Assists the FOAs with data monitoring— FOAs monitor real time data looking for combinations of telemetry parameter values, trends, and other indications that may signify a problem with the satellite or its instruments. The expert systems created with GenSAA will assist the FOAs with the tedious task of data monitoring and allow them to focus on other, higher-level responsibilities during real-time contacts with the satellite. This, in turn, will likely result in more efficient and effective operations.

- Reduces development time and effort; allows quick and accurate response to necessary modifications— The behavior of an orbiting satellite is quite dynamic and occasionally different than anticipated. To quickly create or modify expert systems that can effectively monitor satellites, tools are needed that allow analysts to formulate rule bases easily without the intervention or delay of knowledge engineers and programmers. Several benefits are expected by eliminating these traditional developers. Analysts will be able to create rules quickly in response to unforeseen changes in spacecraft behavior or operational procedures. Also, knowledge translation errors will be reduced or, at least, more easily corrected. Knowledge translation errors are errors which are inadvertently introduced during the process of translating a piece of expert knowledge into rule form.

- Serves as a training tool— In addition to assisting the FOAs with real-time spacecraft operations,

GenSAA will be useful as a training tool in two ways. First, by utilizing the playback utilities provided by TPOCC, analysts will be able to replay a previous spacecraft communications event. Thus, a student analyst can observe how the expert system handles a specific problem scenario. Exercises like this will provide a realistic, hands-on environment for training FOAs in a safe, off-line mode. Second, experience from previous expert system projects indicates that the development of rules used in an expert system is a beneficial mental training exercise for the FOA. When FOAs create rules themselves, they must consider alternatives more closely and may therefore develop a deeper understanding of the problem domain. This approach may enable more effective fault isolation methods to be identified.

- Protects against loss of expertise— Another benefit of automating fault-isolation tasks with rule-based systems is that the resulting rule base serves as accurate documentation of the fault-isolation method. The rule base can be studied by student analysts to learn about fault-isolation techniques. Even more importantly, mission operations can be better protected against the effects of personnel turnovers. POCC expert systems that capture fault-isolation knowledge preserve expertise from mission to mission and mitigate the impact of the loss of experienced FOAs.

Applicability to Industrial Systems

Although initially developed to support real-time data monitoring in satellite control centers, GenSAA will support the rapid construction of highly graphical expert systems for a variety of applications throughout industry. The Rule Builder and User Interface Builder of the GenSAA WorkBench facilitates the development of a knowledge base integrated with a graphical user interface complete with multiple windows, user input graphical objects, and data-driven graphical objects such as meters, gauges, and strip charts. Using these two WorkBench tools, for example, an organization could easily create expert systems to monitor traffic on a computer network or watch over an industrial manufacturing process, searching for problems and providing decision support or corrective action if a problem is detected.

For more complex systems, GenSAA's publish-and-subscribe method of information sharing enables multiple GenSAA expert systems to share knowledge (configuration values, system status, problem diagnoses, data analysis results) for distributed, hierarchical problem solving. For our application in the satellite control center environment, a number of individual GenSAA expert systems are planned to monitor and diagnose the subsystems onboard the spacecraft and within the ground system. These expert systems will publish their results to a "master" expert system which will monitor knowledge from a number of expert systems to provide a high level view and to isolate problems that exist across subsystem boundaries. This approach makes GenSAA applicable to a wealth of commercial, distributed systems such as on-line monitoring and diagnosis of telecommunication switching systems or realtime load control of power distribution networks.

To receive full advantage of the programming-free approach of the GenSAA WorkBench, the third component of the GenSAA WorkBench, the Data Manager, requires minor modifications to support drag-and-drop capabilities with the other WorkBench tools. The following section briefly describes the approach for customizing GenSAA to support other domains.

Integrating GenSAA Into Other Environments

Even without modifications, GenSAA will readily support the development of highly graphical rule based expert systems. However, in order to receive the "programming-free" benefit that this toolset provides, two steps must be taken: 1) the domain data must be formatted to allow the Data Manager to display it and thereby facilitate the drag-and-drop interoperability with the other WorkBench tools, and 2) the Data Interface Subsystem of the GenSAA Runtime Framework must be configured to manage the stream of the data selected with the Data Manager. There are basically two approaches for adapting the Runtime Framework for a new environment:

- Modify the GenSAA Data Interface – In this approach, the Data Interface subsystem of the GenSAA runtime environment is modified to accommodate the existing interface of the data source. The advantage to this approach is that the existing data source remains unchanged. However, the disadvantage is that the new user must modify unfamiliar code (GenSAA) and re-implement these modifications for any subsequent GenSAA releases.

- Create a custom Data Server– Perhaps a better approach for integrating GenSAA with the new data source is to create an intermediary process that functions as a TPOCC Data Server from the perspective of

the GenSAA Data Interface. This process would receive all data requests from GenSAA and forward all data from the data source utilizing the standard TPOCC interface used by GenSAA. Several advantages would result: the group performing the integration does not have to modify foreign (GenSAA) code, updates to the GenSAA tool will not require re-implementation of the customized portions, and conformance to the original GenSAA Data Interface is maintained. The primary disadvantage is the performance penalty that may result from the extra processing in the intermediary process.

Although the modifications necessary to adapt GenSAA to a new environment may initially sound like too much effort, our experience has demonstrated that it is well worth the investment; if multiple expert systems are to be developed, the time spent customizing the front end of GenSAA is easily less than the effort that would otherwise be necessary to integrate each expert system with the corresponding data source. By employing object-oriented design techniques in GenSAA, the modification is simplified and isolated to specific objects thereby preventing the inadvertent corruption of other GenSAA elements that do not require modification.

Conclusion

Detecting satellite anomalies is a challenging task that is becoming more difficult as spacecraft become more complex, the number of sensor points multiplies, and data rates increase. As demonstrated by the CLEAR System, fault-isolation expert systems can help human analysts monitor the flood of data. Expert systems can accurately monitor hundreds of real-time telemetry parameters and isolate discrepancies and anomalies the instant they can be detected. They can alert the analysts while providing advice on how to correct problems swiftly and effectively. Unfortunately, development of these systems is often time consuming and costly; moreover, they usually cannot be reused for other missions.

Consequently, GenSAA is being developed for use by the FOAs who work in satellite control centers. GenSAA is designed to enable fault-isolation expert systems to be developed quickly and easily, and without the delay or costs of knowledge engineers and programmers. By facilitating the reuse of expert system elements from mission to mission, GenSAA will reduce development costs, preserve expertise between missions and during periods of personnel turnover, and provide more effective spacecraft monitoring capabilities on future missions. In the commercial industry, similar benefits can be realized with expert systems, and, although GenSAA was originally developed to assist with spacecraft monitoring, it naturally supports the rapid development and deployment of graphical intelligent monitoring systems in a wide range industrial applications.

References

1. Bielawski, L., & Lewand, R. (1991). *Intelligent Systems Design: Integrating Expert Systems, Hypermedia, and Database Technologies*. New York, New York: John Wiley & Sons.
2. Hughes, P.M. (1989, October). *Integrating Expert Systems into an Operational Environment*. AIAA Computers in Aerospace VII Conference, Monterey, California.
3. Hughes, P.M., & Luczak, E. (1991, October). *GenSAA: Advancing Spacecraft Monitoring with Expert Systems*. The American Institute of Aeronautics and Astronautics Computers in Aerospace VIII Conference, Baltimore, Maryland.
4. Hughes, P.M., & Luczak, E. (1991, May). *The Generic Spacecraft Analyst Assistant (GenSAA): A Tool for Automating Spacecraft Monitoring with Expert Systems*. 1991 Goddard Conference on Space Applications of Artificial Intelligence, Greenbelt, Maryland.
5. Luczak, E.C., Gopalakrishnan, K., & Zillig, D.J. (1989, May). *REDEX: The Ranging Equipment Diagnostic Expert System*. 1989 Goddard Conference on Space Applications of Artificial Intelligence, Greenbelt, Maryland.
6. Perkins, D., & Truskowski, W. (1990, September). *Launching AI in NASA Ground Systems*. AIAA/ NASA Second International Symposium on Space Information Systems, Pasadena, California.
7. Wall, S.D., & Ledbetter, K.W. (1991). *Design of Mission Operations Systems for Scientific Remote Sensing*. Taylor & Francis, Inc., Bristol, Pennsylvania.

N 93-25590
-52.9-61

150499

p. 10

TARGET: Rapid Capture of Process Knowledge

C.J. Ortiz and H.V. Ly
Software Technology Branch (PT4)
NASA/Johnson Space Center

T. Saito
Computer Sciences Corporation

R.B. Loftin
University of Houston-Downtown and
Software Technology Branch (PT4)
NASA/Johnson Space Center

ABSTRACT

TARGET, Task Analysis/Rule Generation Tool, represents a new breed of tool that blends graphical process flow modeling capabilities with the function of a top-down reporting facility. Since NASA personnel frequently perform tasks that are primarily procedural in nature, TARGET models mission or task procedures and generates hierarchical reports as part of the process capture and analysis effort. Historically, capturing knowledge has proven to be one of the greatest barriers to the development of intelligent systems. Current practice generally requires lengthy interactions between the expert whose knowledge is to be captured and the knowledge engineer whose responsibility is to acquire and represent the expert's knowledge in a useful form. Although much research has been devoted to the development of methodologies and computer software to aid in the capture and representation of some types of knowledge, procedural knowledge has received relatively little attention. In essence, TARGET is one of the first tools of its kind, commercial or institutional that is designed to support this type of knowledge capture undertaking. This paper will describe the design and development of TARGET for the acquisition and representation of procedural knowledge. The strategies employed by TARGET to support use by knowledge engineers, subject matter experts, programmers and managers will be discussed. This discussion includes the method by which the tool employs its graphical user interface to generate a task hierarchy report. Next, the approach to generate production rules for incorporation in and development of a CLIPS based expert system will be elaborated. TARGET also permits experts to visually describe procedural tasks as a common medium for knowledge refinement by the expert community and knowledge engineer making knowledge consensus possible. The paper briefly touches on the verification and validation issues facing the CLIPS rule generation aspects of TARGET. A description of efforts to support TARGET's interoperability issues on PCs, Macintoshes and UNIX workstations concludes the paper. Systems such as TARGET has the potential to profoundly reduce the time, difficulties, and costs of developing knowledge-based systems for the performance of procedural tasks.

INTRODUCTION

The nature of their delivery and implementation methods and styles as well as their ability to extract knowledge characterize systems designed to aid knowledge acquisition. Various authoring tools have evolved to solve the problems associated with the creation of a specific expert system [2]. Historically, developers and researchers directed most knowledge acquisition oriented tool designs toward rating or categorizing problems or knowledge. To use such tools to capture specific knowledge, the developer distinguished between types of knowledge methods or approaches. Although sharing many of the same goals, the existing methodologies are numerous, ranging from frame modeling to case-based reasoning models to repertory-grid rating structures. The various knowledge types addressed by these systems—from semantic or taxonomic to declarative to procedural—affect the design and performance decisions of researchers and developers [5]. Knowledge representation, including frames, objects, rules, and decision trees, captures and executes expertise. At this point, most would agree that no one tool accommodates all the cognitive styles needed to gather the information or knowledge necessary for the creation of an expert system in one contiguous process. Clearly, viable standards have yet to be fully established and accepted.

This paper will first discuss the problem of procedural knowledge acquisition and pertinent related issues. Next, the focus will shift to TARGET and its functions and features. Plans for TARGET and trends in knowledge acquisition for the future will conclude this paper.

THE TARGET METHOD

Theoretical Considerations

Procedural knowledge acquisition through task analysis is a reasonable candidate for graphical representation modes. Decomposing a complex set of steps that make up a specific mission or task requires cognitive visualization and the ability to formulate and reformulate the decomposition of those steps or actions. The specific heuristic procedures that most subject matter experts (SMEs) employ share certain levels of organization and recall. The path in which a procedure evolves starts with specific agendas and goals [4]. The last or final action of reaching or satisfying those actual goals ends the procedure. On the other hand, any actions that would restart a process (i.e., a loop) would occur before the goal oriented or last action. Options or implied decisions during a task that direct the expert along alternative paths may or may not affect other performances of the same task. In cases where the processes offer one or more options to complete a task, the process diverges into as many paths as necessary to meet the optional requirements.

The expression of task knowledge could be further divided into strategic and serial styles. To describe specific processes within a domain, a facility should allow the user to express strategic or serial styles of a process. Linear actions pose structural relationships that contrast those that reflect technique or style. Task knowledge manifests itself more often in linear arrangement where pre and post conditions of an event directly couples with the completion of a specific action.

Strategic knowledge combines linear relations within task hierarchies with deterministic actions. It is knowledge employed by an actor in deciding what succeeding action or actions are executable. These actions have consequences external to that actor. The actor's strategy is then elicited from the composition of such actions and their consequences.

This view of knowledge acquisition asserts the construction of abstract descriptions of the tasks themselves, defining the appropriate methods for automating the problem domain. It also serves to apply these methods to domains conforming to the task description [6]. This approach assumes that the knowledge acquisition activity takes place relative to task styles defined *a priori* for the domain. The abstract description of the process or processes generally encodes itself into the task performance environment or actual application. In addition, this representation of the process method serves to constrain the construction of task or domain models to automate the task analyzer for the target application.

The trend toward organizational knowledge has been due in part to a confluence of views between the AI and database communities. This knowledge view is based on the premise that there exists a single knowledge source that may be operated on simultaneously by many different and diverse agents and reasoning engines. Because there is, within reason, a separation of task knowledge from the knowledge of problem-solving tasks that utilize it, the knowledge acquisition problem should be examined within a framework that assumes and supports independence (and to some degree, incognizance) of the abstract problem solving methods and tasks to which the organizational knowledge will be directed.

Role of User Interface and Feedback Mechanism

TARGET exemplifies a genre of knowledge tools that employs directed graph and task decomposition techniques to elicit information for representation in a mode compatible with other development environments like expert system shells, programming or host languages, etc. Task expressive knowledge base architectures are evolving to support common problem tasks or methods at understandable, and subsequently, useful levels of abstraction [3,7]. Task distinct structures could be utilized within software environments that specialize general representations and methods to particular classes of tasks, such as fault detection or medical diagnosis. By abstracting the common characteristics of a class of task, the architecture minimizes redundancy of design as well as elicitation effort. Another supporting determinant would be the insulation of the knowledge engineer and expert from distracting rigors of implementation, permitting concentration on more domain specific issues such as encoding, testing, and refinement of the knowledge base.

TARGET attempts to reinforce the fragile balance between ease of use and design complexity/intricacy of the interface environment. Although not possessing the "bells and whistles" of more sophisticated systems like Aquinas, Protégé and similar tools, TARGET provides enough knowledge modeling (procedural and declarative) support to allow the SME (subject matter expert) or knowledge engineer to build a moderately elaborate knowledge base without sacrificing the attractiveness of its user interface [10]. The intent of this design strategy is for the user or users to employ TARGET's windowed environment to accomplish decomposition of a complex process.

Ultimately, the user, knowledge engineer or SME is responsible for fidelity of the knowledge base before its representation in or transfer to other applications. TARGET's report facilities offer some assistance in this quality checking process. To provide moderately high-level feedback to the knowledge engineer and SME, TARGET can generate the following reports: 1) Task hierarchy -- Sequential or hierarchical account of tasks and 2) User Note Pad - Notes on conditions, states or other user-supplied details.

TARGET supports most cognitive phases of knowledge acquisition with its network approach to knowledge representation. TARGET provides a reasonably comprehensive mechanism for generating simple representations at the very first knowledge acquisition session. Subsequent sessions embellish already elicited knowledge or create new or modified versions of the knowledge base. The user interface gives the user freedom to generate as complex a hierarchy or schema of knowledge as necessary. However, the disadvantage to such freedom is the ability to create a completely abstract knowledge base with relatively few standards for input. Some guiding controls from the TARGET interface could provide structure to the knowledge acquisition process and greatly enhance the ability of the user to create a useful knowledge base.

FEATURES AND FUNCTIONS

Current Interface Capabilities

This section will describe overall features and functionality that TARGET offers as part of the knowledge modeling activity. Portability, a major issue concerning TARGET's delivery strategy, is discussed later in this paper.

TARGET provides the user a dialogue box for task description and selection of a task type. As the user enters text description in free form, the user may be as descriptive as possible. If more explanation of a task is required, the user can attach a notebook entry for additional information. Tasks with notebook entries attached all have a small musical note located in their lower right hand corner.

The arrows linking the tasks describe a process and a flow of control between the tasks. TARGET scans the network for connectivity errors that users may have inadvertently introduced into the system. The system will jump to the level where the error was first encountered and display the responsible task. A descriptive error message comes from a pop-up dialogue box.

Upon completion of the network and satisfaction of all connections, the system generates a work breakdown or a task hierarchy report. This process assigns sequence numbers to each task and will display the final report to the screen, printer or save it as an ASCII text file.

Figure 1 illustrates the basic design interface components of TARGET. Located along the left hand side of the interface, the 'toolbox' contains all the icons used to create and manipulate objects when creating a network. A task-box icon (first button) creates a new task or allows the user to edit an existing one. Next, a linking button links tasks by creating an arc from one task to another and removes links between tasks. The moving van icon selects tasks to be moved from one location to another. Users can copy tasks by selecting the copy button. The trash can icon selects tasks from the current layer for deletion or removal.

As the number of tasks increases, so does the complexity of the network. Users can utilize the Zoom facility from the main menu to move the field of view to a wider angle, resulting in the tasks appearing smaller on the screen. The magnifying glass in the toolbox selects and magnifies a given task bringing it into full view. The tree icon shows the overall hierarchy of the network allowing users to select and quickly jump to any level in the network. The final toolbox button provides context sensitive help. By selecting the question icon, users can gain access to the on-line hypertext reference manual containing topics on every icon, object, or menu item in the user interface.

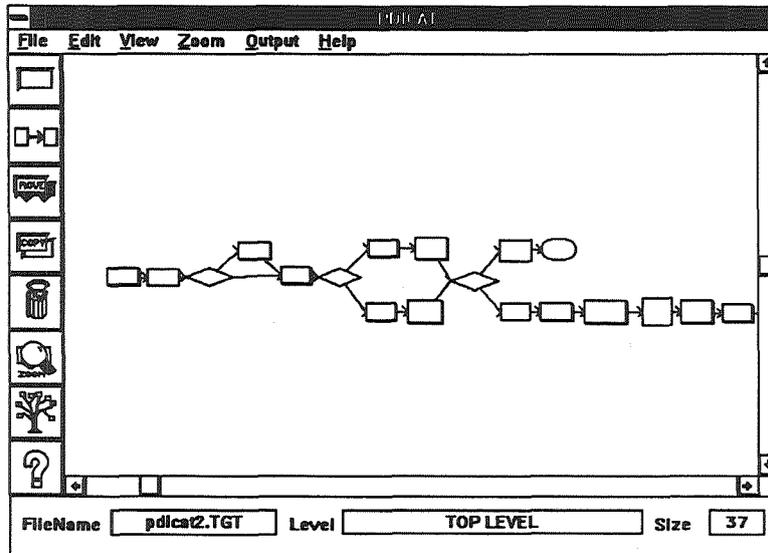


Figure 1. TARGET Interface

The developers incorporated an on-line version of the TARGET reference manual to enhance the interface and knowledge acquisition session. Users may click the help icon to select any object, menu item, or task for help directly from the on-line manual. The hypertext help engine generates context sensitive help on all Windows-based platforms.

Version 1.5 Features

The user interface has undergone several changes as TARGET has gone from version 1.0 to 1.5. The icon bar for task manipulation has been relocated to a ribbon configuration along the top of the screen. The developers improved task manipulation with the ability to select individual as well as groups of tasks.

TARGET no longer limits task number or size displayed at any one time. TARGET now incorporates notebook and other information directly into the task hierarchy report with a better numbering and printing algorithm. Along with textual reports, TARGET provides the ability to print the graphical information. The internal design allows the ability to reuse previously created procedures and link to external files along with a portable version of the on-line help reference manual.

PROCEDURAL EXPRESSION

Graphical Task Algorithm

A user describes task information required as input and the task is then characterized by the information it produces as output. Let $C = \{c1, \dots, cn\}$ be a set of n tasks, where n is some integer. Let $P = \{p1, \dots, pm\}$ be a set of m parameters. Let $V = \{v1, \dots, vq\}$ be a set of locations with q as small integer values. Let d be a function that translates parameters P to the values V , $d:P \rightarrow V$ (value assignment). For example, parameters P are sets of assertions within a level. The value s is sequential value of each task entry where d is the coordinate value. The values V are defined as relational values, in which case $d(p)$, the value of parameter p , would be the sequence value of assertion p . Since the number of values q can be more than 2, this allows for specification of incomplete knowledge. This does not entail any particular loss of generality as long as V can accommodate the range of values for every parameter and as long as each parameter can be associated with a different interpretation of values. Within TARGET for instance, redundant occurrences of C or an action kernel would be controlled by a four tuple: $\langle P, V, C, s \rangle$. In this case, s would be the function from D to C , $s:D \rightarrow C$, where D is the set of value assignments of the parameters P to the values V . The resulting output would still reflect the distinct permutations of actions based on input parameters of the simple implementation of s .

Where directed graphs are concerned, the adjacency matrices may be used to describe the links or arcs between nodes. The use of rule representation requires the specification of the adjacency matrix that connects the appropriate nodes.

As a result, the matrix bounded by the number of unique clauses and rules will represent the graphical configuration [10]. TARGET provides link and coordinate data that directly correspond with pre and post conditional relationships. At the outset, TARGET's intermediate knowledge representation, in the form of its task hierarchy, establishes the tone for ultimate translation into the knowledge base. To further elaborate, a rule set will consist of n rules. If an individual rule is represented as ri , then let the antecedents and consequents of this rule be represented by $AC(ri)$ and $CC(ri)$ respectively. Individual clauses are designated as $AC_j(ri)$ and $CC_k(ri)$. Consequently, a rule of the following form would involve nodes $AC_1(ri), \dots, AC_m(ri), CC_1(ri), \dots, CC_n(ri)$, and ri :

$$AC_1(ri) \& \dots \& AC_m(ri) \rightarrow CC_1(ri) \& \dots \& CC_n(ri)$$

This would set to 1 the values in the adjacency matrix corresponding to nodes:

$$(AC_1(ri);ri), \dots, (AC_m(ri);ri), \dots, (ri;CC_1(ri)), \dots, (ri;CC_n(ri)).$$

Task Types within the TARGET Domain

TARGET classifies each task within its dialog box in the following types:

Task Type	Function	Shape
1) Required	Must be completed when first encountered	Lightly colored rectangle (Figure 4)
2) Optional	Not required; If performed, must be executed before a specified action	Shaded or gray rectangle
3) Goal	Ends a sequence	Rectangular box with a thick border (Figure 5)
4) Control	Provides logical jumps to specified tasks	Oval shape (Figure 6)
5) Decision	Actions based on a set of valid responses	Diamond shape with labeled arcs radiating out

TARGET also provides parallel and parent relationships as basic building blocks to describe a process flow. Parallel tasks are tasks that have multiple arcs radiating from them (Figure 8). Parallel and decision tasks are the only tasks that have more than one path radiating out from them. Subtasks consist of required or optional tasks decomposed into several smaller layers. The task boxes appear with the required or optional qualities with an additional shadowed form representing depth.

TARGET Rule Conversion Mechanism

TARGET, as originally conceived, will produce rules for incorporation into an ICAT (Intelligent Computer-Aided Training) expert system architecture. The procedural rules interact with, and are controlled by, other CLIPS systems using a blackboard architecture. TARGET takes the relationships between actions and organizes them into their antecedent and consequent positions. The format in figure 2 illustrates the most basic rule construct in pseudo code. Next, within CLIPS syntax, generated rules will follow the ICAT oriented format in Figure 3. The logical sequence of actions embeds in the step triggers fact assertions activating the following rule.

<pre>(defrule name ?step <- (previous task has been completed) => (retract the previous task from the fact list) (do zero or more functions (printout, assert, etc.)) (assert a fact that this task has completed))</pre>	<pre>(defrule control_rule ?step <- (next_step ?number) => (retract ?step) (assert (step ?number)))</pre>
Figure 2. Rule Template in English	Figure 3. Simple Control Rule Format

For simple sequential relationships (Figure 4), task A executes before task B. Figure 4 shows the dependencies of the task B upon the successful function and performance of task A. A rule derived from a goal task (Figure 5) will retract the previous control fact and process functions, but will not assert a control fact, thus ending a process.

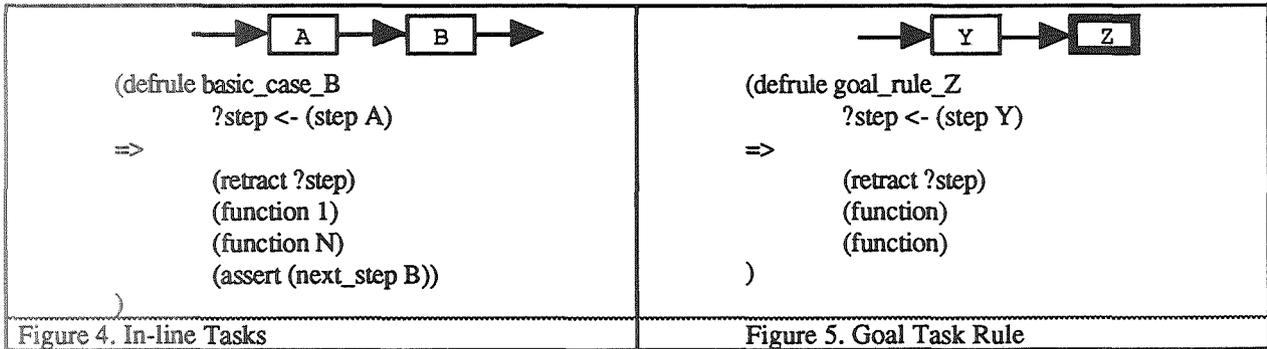


Figure 6 shows that control task rules will fire only after completion of the previous task but will not process a function.

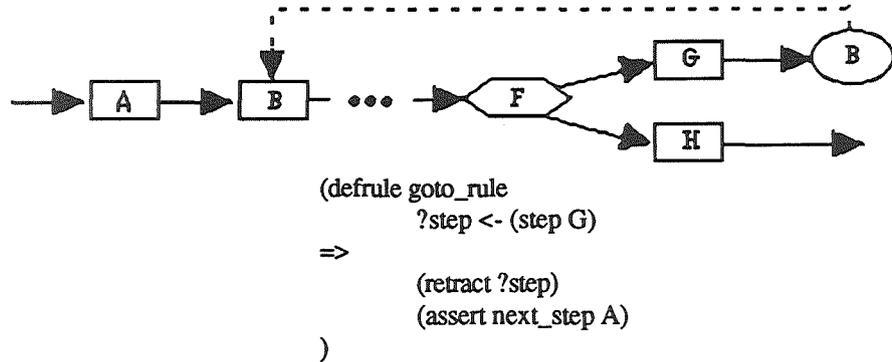
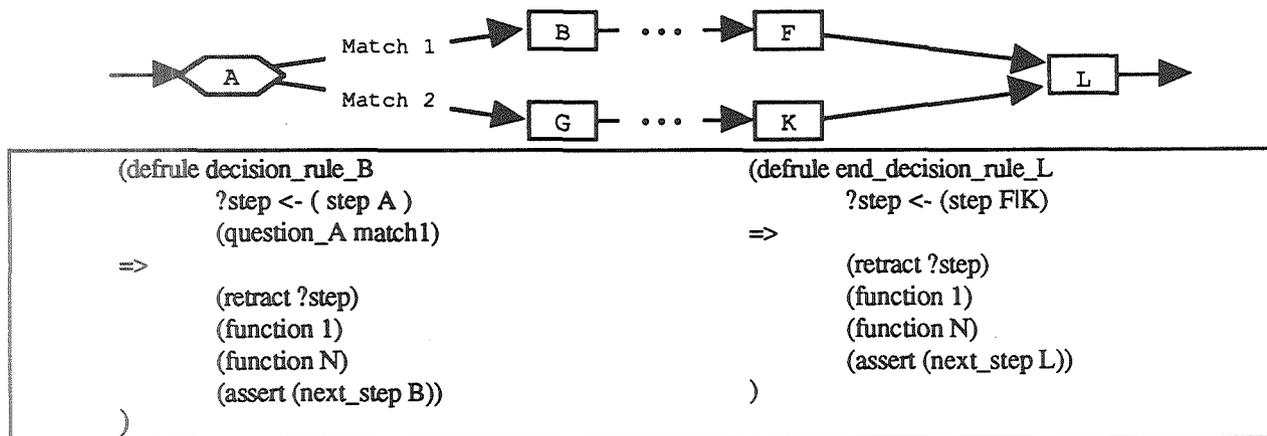


Figure 6. Control Task Rule

TARGET can generate more complex rules involving decisions and branching. In the following example (Figure 7) decision A has two possible answers represented by Match 1 and Match 2. The outcome of the decision activates one path (B or G). As the two paths converge, task L verifies successful completion of either task F or K.



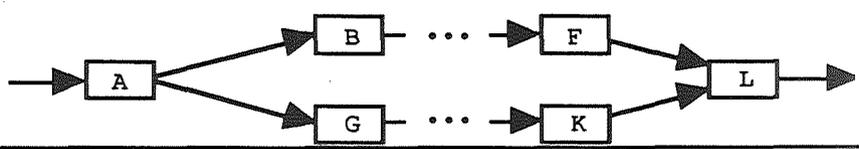
```

(defrule decision_rule_G
  ?step <- ( step A )
  (question_A match2 )
  =>
    (retract ?step)
    (function 1)
    (function N)
    (assert (next_step G))
)

```

Figure 7. Parallel Alternative Paths

Although CLIPS rules generated for parallel tasks are similar to those generated by decisions, they differ in the following qualities (Figure 8): 1) once Task A fires, rules B and G are both activated and 2) the control fact for rules B and G from the previous step is not retracted from the fact list. Finally, the rule associated with task L must not fire until tasks F and K have completed.



```

(defrule parallel_rule_B
  ( step A )
  =>
    (function 1)
    (function N)
    (assert (next_step B))
  )
)
(defrule parallel_rule_G
  ( step A )
  =>
    (function 1)
    (function N)
    (assert (next_step G))
  )
)
(defrule end_decision_rule_L
  ?step1 <- (step F)
  ?step2 <- (step K)
  ?start <- (step A)
  =>
    (retract ?step1 ?step2 ?start)
    (function 1)
    (function N)
    (assert (next_step L))
  )
)

```

Figure 8. Parallel Simultaneous Paths

SEAMLESS PLATFORM DELIVERY

Portability became a crucial aspect of TARGET's delivery method after the Version 1.0 release.

Multi-platform Delivery

To support current and potential users, TARGET must be able to run on a wide variety of machines including PCs, Macintoshes, and Unix workstations. To reduce the cost in creating independent user interfaces for all platforms within their native development environments, the TARGET development team has employed XVT (Extended Virtual Tool Kit developed by XVT, Incorporated) to implement the user interface for TARGET. XVT is a set of graphical functions that resides in the native development environment of each system. XVT allows users to develop and maintain a single C or C++ application that runs on the twenty-six different systems that XVT supports. TARGET/XVT currently runs, in test mode, on Macintosh (MacOS), IBM PC compatible (Microsoft Windows) and Sun (Motif) platforms. XVT should also permit TARGET to run in MS/DOS, Open Look and OS/2 environments as well.

Originally developed in the Microsoft Windows 3.1, Software Development Kit environment, TARGET's kernel code is nested within the user interface code. For the sake of maintenance and flexibility, developers ported the original user interface code to XVT and separated the kernel from it. This configuration protects the kernel from any XVT software upgrade. For XVT upgrades, only the user interface code requires modification. Similarly, the developers are still able to modify the kernel independently when necessary.

User Benefits

TARGET's goal is to provide a user friendly interface in an environment with which an end user is already familiar. Whether users are comfortable with UNIX-based X-Windows, P.C. or Macintosh platform environments, TARGET can function in all three. Knowledge files generated by TARGET are also portable. Thus, procedures generated by a Macintosh can be parsed and displayed on PC or Unix workstations.

Manipulation of tasks within the interface operates using a one-button mouse mechanism. Although most X-Windows workstations and PCs use multiple-button mouse devices, most Macintosh computers use single button varieties. The user selects, links, moves, copies and deletes tasks through a single mouse click. On the other hand, users maneuver through the task network by double-clicking actions.

Another goal was to eliminate the use of color to convey information. Users with black-and-white output devices such as laser printers need not worry about loss of color information. Within TARGET, shapes and their connectivity convey information about a task type and structure. For users with monochrome monitors, related hardware upgrades become unnecessary.

IMPORTANT ENHANCEMENTS

Notebook facility

Several notebooks can be attached to a specific task. Ability to store audio, video, and still photo information within a notebook is currently being investigated. The developers must address the issue of portability before integrating these advances. The TARGET file structure must be portable among the various workstations. Consequently, audio, video, and pictorial information must be of a standard format. Still photos attached to a task will most likely be stored in the Graphics Interchange Format or GIF. Commercial GIF viewers exist on all standard systems to render images to the screen. Research is ongoing to find or develop similar standards for video and audio information.

Other Modeling Functions

Problem-solving methods can be regarded as knowledge that establishes and controls sequences of actions required to perform tasks or processes. This control knowledge defines the order in which subtasks and subfunctions are resolved to perform and complete an overall task. Although problem resolution comprises procedural mechanisms, other parallel activities may be in operation. Diagnosis or decision making actions may also be in progress [8].

Within TARGET, the kinds of domain-specific knowledge that are applicable within each kernel action help define the problem criteria. Ultimately, the problem solving method identifies and, eventually, classifies the domain knowledge. The granularity of TARGET's problem-solving method hinges on the knowledge characterized by the SME's role within an application without further control or meta-knowledge. It would make the different roles of SME's knowledge bases within a design evaluation task explicit. Ways could also be suggested to organize knowledge base activities according to such knowledge perspectives.

TARGET will provide a fault detection, isolation and resolution (FDIR) facility to compliment its nominal path modeling capabilities (Figure 9). The FDIR component is suited to describing problem solving methods between actions whose antecedent conditions were not instantiated. From that point an implied looping relationship emanates until the problem or issue is resolved. For example, turning the key in an automobile ignition may either produce the desired engine turnover or the process of troubleshooting to determine why the car would not start. In either case, the antecedent provides impetus for the looping mechanism until the consequent is actualized. Then, the operation can proceed to the next step.

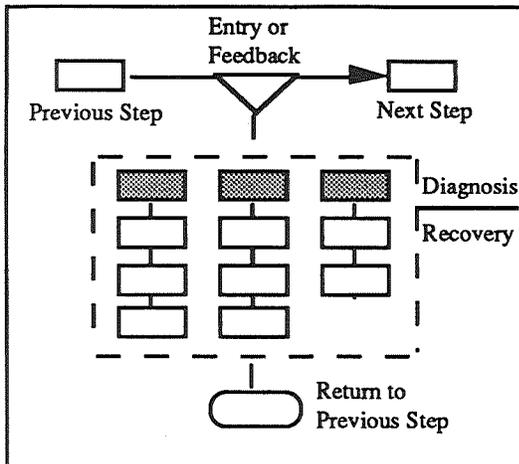


Figure 9. Diagnosis and Recovery Mechanism

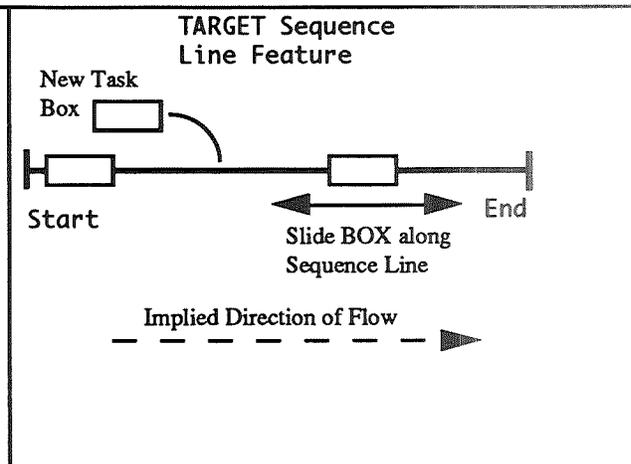


Figure 10. Sequence Line Modeling

A new sequential modeling feature (Figure 10) will be added to TARGET's interface addressing the design problem of layout adjustment while minimizing linking keystrokes for the user. As a result, this will reduce task modeling efforts by automatically creating arcs between tasks. All that the users will see is a sequence line at every logical level to which they will attach and order tasks. The next generation of TARGET will possess sequential intuition for quicker prototyping of procedural relationships.

CONCLUSION

Potential Applications

TARGET offers the potential of application to other areas based on its hierarchical reporting mechanism and graphical design. Other than the ICAT environment applications, potential users of Total Quality Management (TQM) could benefit from TARGET's ability to study, streamline or expand processes through the graphical interface. The CASE world already uses box-flow techniques to model procedural language development where TARGET could affect several issues, including code generation, reverse code engineering and browsing.

Knowledge Capture Technology

TARGET could significantly impact the development of various ICAT systems as well as the development of other intelligent systems. For any procedural knowledge acquisition task, it can enhance the ability of the expert to visualize and organize a task or process. Procedural visualization of this type will become more popular as more tools with organizational diagnosis capabilities evolve [1].

As computer hardware power increases, more latitude in presentation methods will be available. Visual conception and communication of abstract information will become more common. The strategic fusion of graphical display (bit-map, meta-graphic, etc.) and graphical input device (mouse, light-pen, trackball, etc.) technologies will facilitate visual as well as textual representation of knowledge [9]. Drawing tools already allow the user to produce and manipulate complex graphics. The role of these tools can also combine with organizational algorithms to create more intelligent diagrams, flow charts, and interactive decision trees. With users becoming more adept at employing systems with pictorial modeling capabilities and computers better able to support complex graphical interfaces, the "TARGET" mode of procedural knowledge acquisition will become more widely accepted.

As knowledge acquisition evolves as a discipline within artificial intelligence, more tools to assist in the knowledge acquisition process will also become available in useful forms. TARGET, and tools like it, will be employed within their own "niche" and will also be integrated with other methodologies in the future [11]. Although TARGET currently models the sequence within the task hierarchy structure for rule induction, additional efforts will be devoted to encapsulating additional knowledge into the steps within a network. In particular, address issues such as gathering artifact data, selected action rationale, and interactive verification and validation of rules will be addressed.

REFERENCES

- [1] Akscyn, R. M., McCracken, D. L. and Yoder, E. A. (1988). "KMS: A Distributed Hypermedia System for Managing Knowledge in Organizations", *Communications of the ACM*, July, 31 (7), 820-834.
- [2] Boose, J. H. (1989). A survey of knowledge acquisition techniques and tools. *Knowledge Acquisition*, March, 1 (1), 3-37.
- [3] Clancy, W.J. (1985). Heuristic classification, *Artificial Intelligence*, 27, 291-349
- [4] de Kleer, J., Doyle, J., Steele, G. L., Jr. and Sussman, G. J. (1985). AMORD: Explicit Control of Reasoning. in *Readings in Knowledge Representation*, Brachman, R. J. & Levesque, H. J., Eds., Los Altos, CA: Morgan-Kaufmann Publishers, Inc., 345-355.
- [5] Gaines, B. R. (1988). An overview of knowledge-acquisition and transfer. in *Knowledge Acquisition for Knowledge-Based Systems*, Gaines, B. R & Boose, J. H., Eds., *Knowledge-Based Systems, Vol.1*, New York: Academic Press, 3-22.
- [6] Gruber, T. R. (1989). *The Acquisition of Strategic Knowledge*, Academic Press, Inc., Harcourt Brace Jovanovich.
- [7] Gruber, T.R., and Cohen, P.R. (1987). Design for acquisition: Principles of knowledge system design to facilitate knowledge acquisition, *International Journal of Man-Machine Studies*, 26(2), 143-159.
- [8] Kitto, C.M., and Boose, J.H. (1988). Heuristics for expertise transfer: An implementation of a dialog manager for knowledge acquisition, *Knowledge Acquisition Tools for Expert Systems, Knowledge Based Systems*, Academic Press Limited, 2, 175-194.
- [9] Messinger, E. B., Rowe, L. A. and Henry, R. R. (1991). A divide-and-conquer algorithm for the layout of large directed graphs. *IEEE Transactions on Systems, Man, and Cybernetics*, 21 (1), 1-11.
- [10] Musen, M. A., Fagan, L. M. and Shortcliffe, E. H. (1986). Graphical specification of procedural knowledge for an expert system. *Proceedings of the 1986 IEEE Computer Society Workshop on Visual Languages*, Dallas, Texas, 167-178.
- [11] Saito, T., and Loftin, R. B. (1990). Supplemental knowledge acquisition through external product interface for CLIPS. *First CLIPS Conference Proceedings*, NASA Conference Publication 10049, 1, 174-179.
- [12] Shaw, M.L.G. and Gaines, B.R. (1987). Kitten: knowledge initiation and transfer tools for experts and novices. *International Journal of Man-Machine Studies*, 27, 251-280.

N 9 3 - 2 5 5 9 1
530-6/1
150500
P 10

Tree Classification Software

Wray Buntine, RIACS
NASA Ames Research Center
Mail Stop 269-2
Moffet Field, CA 94035

ABSTRACT

This paper introduces the IND Tree Package to prospective users. IND does supervised learning using classification trees. This learning task is a basic tool used in the development of diagnosis, monitoring and expert systems. IND was developed as part of a NASA project to semi-automate the development of data analysis and modelling algorithms using artificial intelligence techniques. IND integrates features from Breiman *et al.*'s CART and Quinlan's C4 with newer Bayesian and minimum encoding methods for growing classification trees and graphs. IND also provides an experimental control suite on top. The newer features give improved probability estimates often required in diagnostic and screening tasks. The package comes with a manual, Unix ``man'' entries, and a guide to tree methods and research. IND is implemented in C under Unix, and has been beta-tested at university and commercial research laboratories in the United States.

DIAGNOSIS AND CLASSIFICATION

A common inference task is where we learn to make a discrete prediction about some case given other details about the case. For instance, in financial credit assessment we wish to decide whether to accept or reject a customer's application for a loan given particular personal information. In monitoring a subsystem of the space shuttle, measurements such as flow rates and temperature are continuously recorded and we need to screen those measurements to decide if the system is in normal or abnormal operation. If the system is in abnormal operation we might further wish to try and predict the type of abnormality present. This prediction task is the basic task of many expert systems, health monitoring systems, diagnostic systems, etc. Furthermore, more complex problems can often be broken down into a sequence of simple prediction problems. For instance, speech understanding, converting the spoken word into written text, is a sequence of prediction tasks about each phoneme.

In medical diagnosis, or diagnosis of equipment subsystems, we need more than just a prediction, we need a careful probabilistic assessment. A simplistic medical example will bring this point home. Suppose your doctor suspects you have a cyst in your abdomen. The options (1 or 2) and outcomes (A or B) give the following set of possibilities: (1A) operates, discovers a cyst, removes it, and you're grateful; (1B) operates, no cyst found, but you're left with the medical bill and a day recovery in hospital; (2A), doesn't operate but the cyst exists and causes medical complications due to lack of treatment; (2B), doesn't operate, no cyst exists. Each case has important implications to you both financially and in quality of life. With a careful probabilistic assessment of the existence of a cyst, you can weigh up the options and decide which option (1 or 2) is the most beneficial to you. For instance, if the medical bill is insignificant compared to the potential medical complications, then you would decide to have the operation even if there was a small chance of having the cyst. If the potential medical complications were insignificant, you would only decide to operate if there was a very high probability of having the cyst. This process of decision analysis requires as input probabilities about the new case in question.

In health monitoring and diagnosis, these probability assessments are needed when the system is being used to screen cases, i.e. the computer systems scans the on-line monitoring data and at certain time points alerts a human expert that a potentially anomalous situation has arisen. Probability assessments such as the "probability of equipment failure" can be used to determine which of the many cases scanned should be forwarded to the human expert for the more costly process of manual inspection.

I will refer to this prediction problem as *classification*, where the aim is to *classify* each new case. One common technique for developing a system to do prediction or probability assessment about new cases is to examine a database of cases, for instance collected historically. Assume that hindsight tells us which is the correct classification for each case in the data base, so for each we know which prediction was optimal. From the data base

we use statistical techniques to "discover" or "learn" how to do the predictions for new unseen cases. This learning technique is represented in Figure 1. The process requires three main forms of input: an *expert* who is able to advise on the problem, help configure the system, etc., a *data base* of correctly classified cases to use in the learning process, and a *model family* from which the learning algorithm is to select a "good" model for doing prediction or probabilistic assessment.

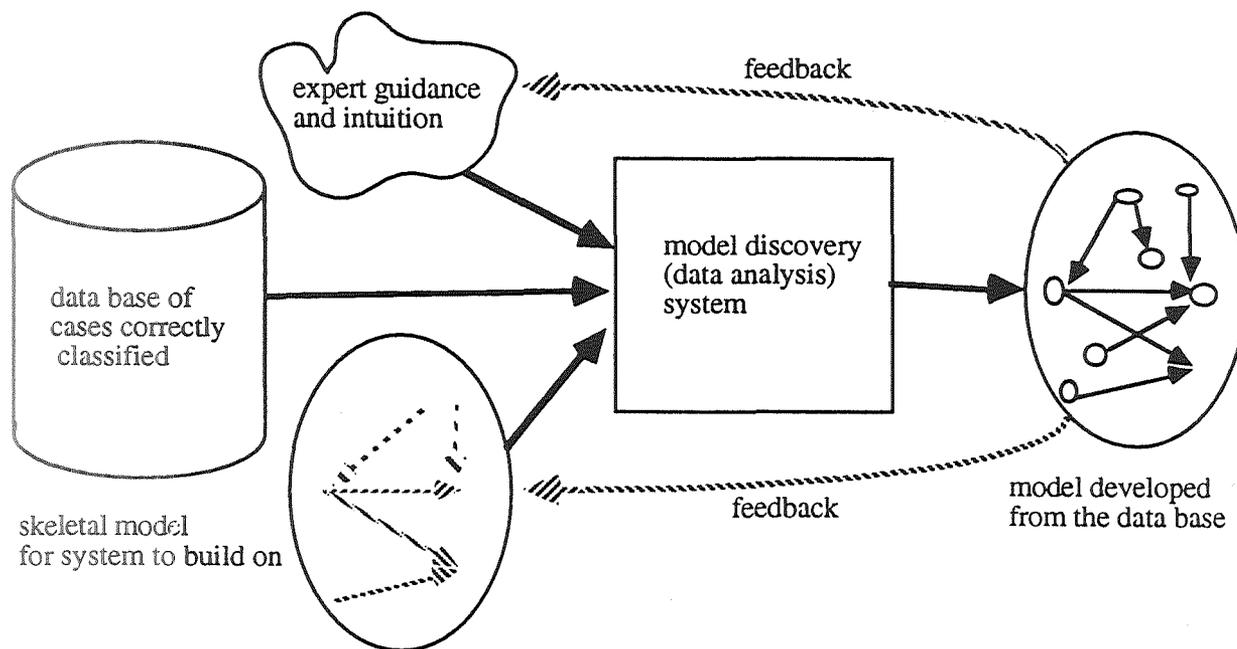


Figure 1. Learning prediction models from data.

This model learning or discovery process is a useful technique in almost every industry, finance, manufacturing, etc., wherever on-line databases are stored and important predictions have to be made on a regular basis about new cases before they enter the data base. Not surprisingly, there are many different fields of science that address this problem as one of their central concerns. In artificial intelligence it is referred to as the classification or induction problem. Techniques include tree and rule learning algorithms of the form I will present in this paper. In statistics it is referred to as the discrimination problem, and common techniques are the linear models used in the finance and banking industry for credit assessment. In pattern recognition it is referred to as supervised learning. In neural networks it is the classification and generalization problem and is routinely investigated using a number of network architectures. These diverse fields are all studying the same problem, "learning to predict", and present a confusing array of methodologies and paradigms for addressing that problem. They differ in the following aspects:

Model family: Which class of models are being used to do prediction? In Figure 1 this corresponds to the "skeletal model". I present classification tree and classification graph model families in this paper.

Statistical philosophy: How is learning to occur? That is, what statistical principles if any are used to develop the central box in Figure 1?

Computational and optimization methods: What are the basic computational methods used in terms of efficiency, optimality, search method, etc.?

Methodological support: What methodology does the analyst use to go about applying the technique to a real problem? For statisticians this is the "consultancy phase" rarely covered in university courses. In artificial intelligence this is the process of "knowledge engineering".

I will refer to the general task of learning how to predict (or estimate probabilities) from data as the classification task. The next section discusses the design of tools for this task. After this, the model family considered in this paper is addressed, and the IND program presented.

DESIGN OF CLASSIFICATION TOOLS

This research is part of a broader effort to semi-automate the development of classification algorithms. The goal of this research is to develop generic tools for learning from data and from partial models of the domain, and to develop the capability to rapidly develop and tailor these learning tools for particular domains given, for instance, specification of the kinds of models that are of interest. When encountering a new application, we sometimes find that off-the-shelf-tools, such as IND, need some modification in order to better suit the task. A good development methodology lets that be done with minimum fuss.

Rather than following a particular field, our group takes a multidisciplinary approach and combine a range of methods required to address the classification task. Our group uses artificial intelligence search techniques for discrete search problems, and standard numerical techniques for continuous problems. We use some of the flexible knowledge representation schemes from artificial intelligence as skeletal models or model families (see Figure 1), and use Bayesian statistical and decision methods for the statistical philosophy underlying our learning algorithms. This methodology allows rapid development of approximately optimal algorithms and so avoids the many pitfalls of *ad hoc* development according to "hunches" and the time-consuming refinement cycle that this entails. This theoretical framework of "statistical philosophy" plus "optimization methods" is important because empirical, *ad hoc* development of algorithms in neural networks and early machine learning has been time consuming and is often plagued by unexplained problems. Empirical validation of our algorithms is also important to check approximations made in interpreting the Bayesian theory. We do this empirical validation by applying the algorithms to a battery of recognized learning problems taken from the literature, or manufactured problems. A summary of our groups development methodology is given in Figure 2. This has led to the development of a number of sophisticated algorithms, one of which was the Autoclass system, show-cased at Technology 2001 by Stutz, Cheeseman, and Taylor at San Jose, December 1991.

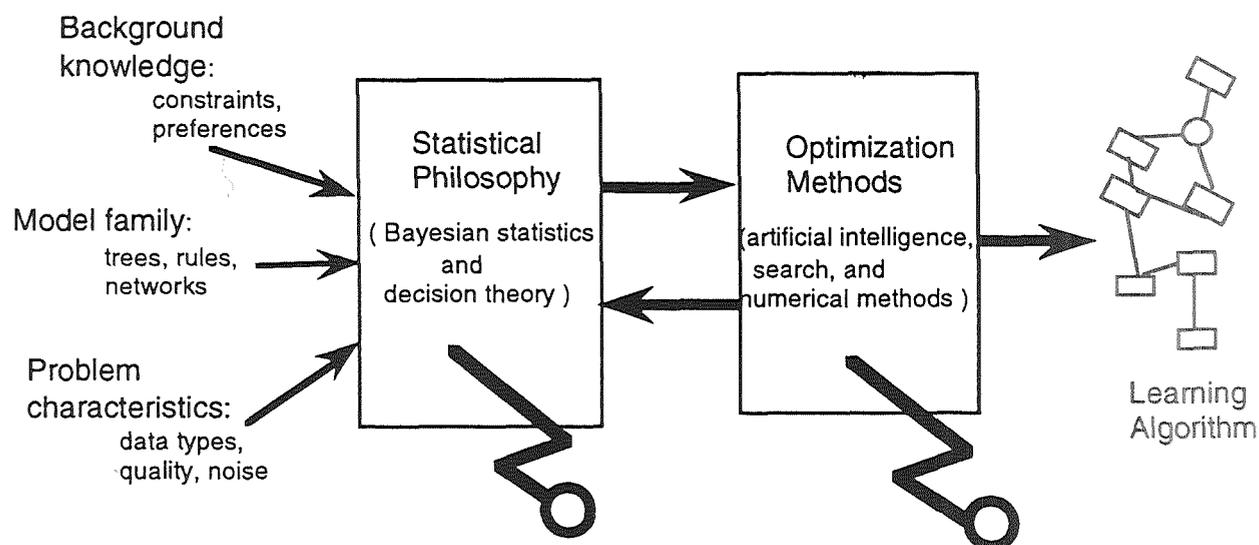


Figure 2. Semi-automatic development of learning algorithms

The justification for Bayesian decision theory, used in the first box in Figure 2, comes from fundamental principles of how uncertain reasoning should be done [1]. The theory applies widely in inference and plausible reasoning and its use is continually expanding in artificial intelligence and neural networks. But there is not a single "Bayesian learning algorithm," as some people mistakenly believe when they learn about the simple Bayesian classifiers developed in pattern recognition. Rather, Bayesian decision theory presents computational guidelines on how learning should be done for many different learning problems, and shows how to tailor methods to particular applications. This means our algorithms can be fine-tuned to the requirements of an individual application. IND has some basic features that allow such tuning.

CLASSIFICATION TREES AND GRAPHS

The IND package described later does prediction using *decision trees* or *decision graphs* and does probability evaluation using *class probability trees* or *graphs*. These are a general form of classification rule that mix discrete and continuous data and are often suited when there is believed to be some form of non-linear structure in the data. A decision tree is shown in Figure 3b. This has the classes *hypo* (hypothyroid) and *not* (not hypothyroid) at the leaves. This tree is for a *two-class classification problem* because there are two different classes that leaf nodes recommend. This tree is processed as follows. Look at the new case you wish to evaluate. Is its value of *TSH* greater than 200? If so take the left branch of the tree and you have reached a leaf. The tree says to predict *hypo*, i.e. hypothyroid. If however the value of *TSH* was less than 200, then take the right branch. Now you have a subtree and you repeat the process. In this case is *Pregnant* true or not? This tree is referred to as a "decision tree" because decisions about class membership are represented at the leaf nodes. Notice that the real valued attributes *TSH* has been incorporated into the tree by making a binary test of the form "attribute < cut-point". Also, the tree need not be binary; if a 4-valued attribute is tested at one of the nodes, then the tree might have 4 branches coming from the node, one for each value.

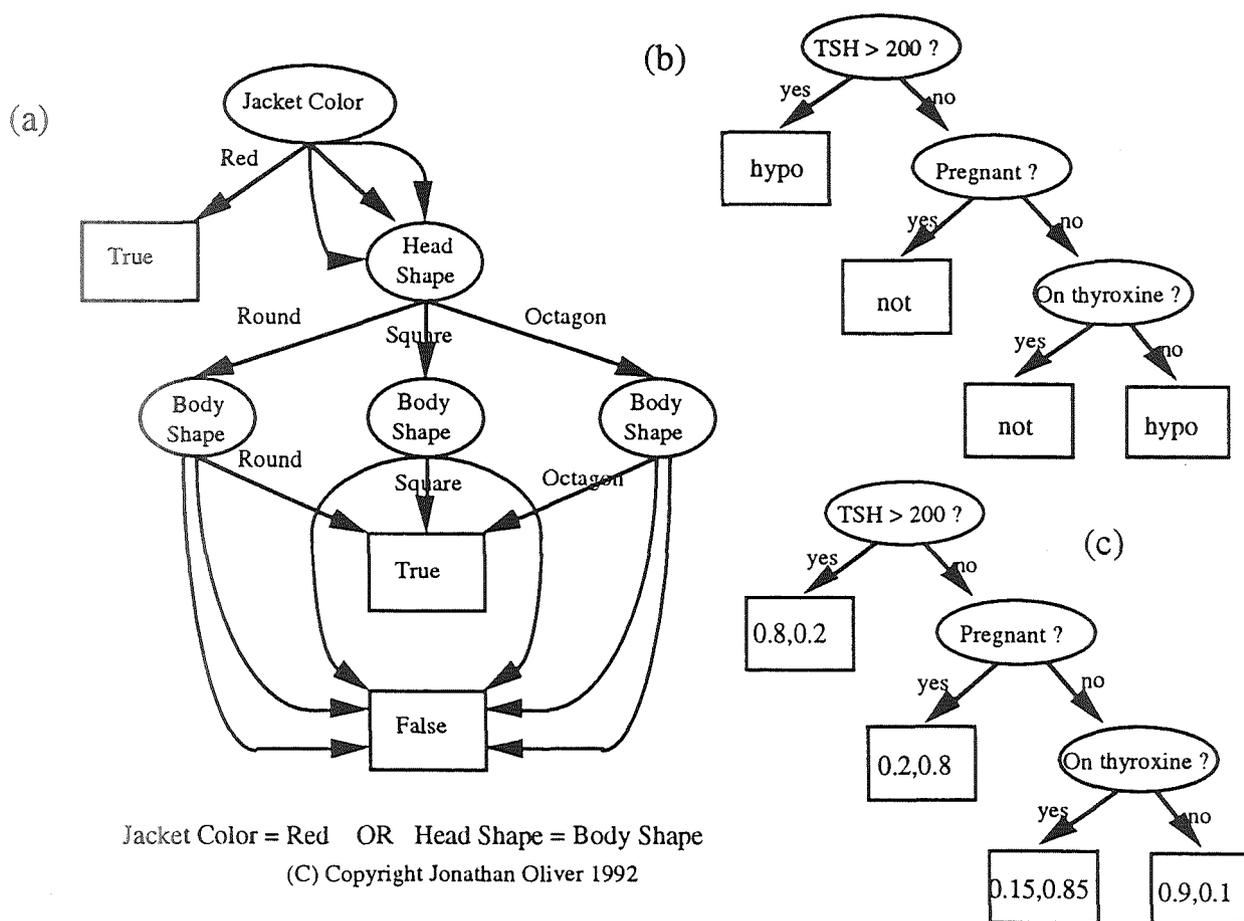


Figure 3. (a) This is a decision graph for the boolean problem given in the figure. Start at the root at trace through the graph to arrive a decision. (b) This is a decision tree for the "hypothyroid" application. (c) This is a class probability tree. Leaf nodes give estimates of class probability.

In typical problems involving noise, class probabilities are usually given at the leaf nodes instead of class decisions, forming a *class probability tree* (where each leaf node has a vector of class probabilities). A corresponding class probability tree is given in Figure 3c. The leaf nodes give predicted probabilities for the two classes. Notice that this tree is a representation for a conditional probability distribution of class given information higher in the tree. This is the statistical interpretation of the tree that Bayesian methods use in developing a learning algorithm.

Methods for learning decision trees and class probability trees have been under development in some form or another for some two decades. The standard technique for building classification trees from data is the so-called recursive partitioning algorithm that forms the basis of systems such as Quinlan's ID3 and C4 [2,3], well known in the machine learning literature, and Breiman, Friedman, Olshen and Stone's CART [4], well known in the applied statistics literature. These methods are largely reimplemented in IND.

A more complex structure is shown in Figure 3a. This is a decision graph, and it is also for a two-class problem. Graphs and trees can also be applied to problems with three or more classes. The graph is processed in exactly the same way as a decision tree, however notice that the graph allows more general connections. This graph represents the boolean function "*jacket-color = red or head-shape = body-shape*". This function would take a complex tree to represent. With graphs we can represent concepts more efficiently. Methods for learning decision graphs and class probability graphs have only recently appeared, and they supersede trees in that they are a more general representation. IND version 3.0 will include these methods, coded up by Jon Oliver [7].

AN OVERVIEW OF THE IND PACKAGE

IND is a suite of C programs and C shell scripts for building tree classifiers and graph classifiers of the kind just described. Currently, several different methods are integrated (CART style; the regression aspect of CART is not implemented, early C4 style, MML/MDL, and Bayesian averaging). Careful checking has been done so that IND reimplements CART and the early C4 fairly faithfully.. The new Bayesian/MML/MDL features can give performance improvement over these in many cases when used appropriately.

IND can be operated in a variety of modes that allow the novice to build trees without too much fuss, and also allow the expert to fine tune the algorithms to particular applications. If you're interested in applying IND to applications, advice is given in the manual on which options to use and how to take into account features of your application and data when configuring your use of IND. If you're interested in running comparative trials or just experimenting with tree software, IND provides extensive experimental control (random partitioning, cross validation) and significance testing. The code for IND is provided (and sometimes even moderately documented) so you can develop your own extensions.

The IND Manual: "An Introduction to IND and Recursive Partitioning" is the best place to start if you are unfamiliar with IND or recursive partitioning. The manual contains an introduction to IND that walks through a few typical sessions, a tutorial on recursive partitioning, a description of IND options, and a fairly complete glossary and bibliography. There is an enormous literature on decision trees and their applications so the manual also contains a brief guide to the literature.

IND has a variety of features including: interactive control of tree building, variable search such as multi-ply look-ahead, missing values and subsetting, handling of utilities and cost functions, prediction of error rates and utilities, a range of priors for the Bayesian methods, printing options, a classifier, etc., a user manual, and a state of the art guide to tree learning research. IND has the look and feel of a typical Unix system and comes with "man" entries. The system has been developed exclusively in a SUN workstation environment under various releases of SunOS UNIX. It compiles under Kernighan and Ritchie C, `cc` and `gcc`, although in future will be converted to ANSI standard C. Various users have ported the system to HP, IBM and other Unix platforms and their changes have been incorporated in the latest release.

In November 1991 the IND Tree Package version 1.0 was prepared and released as a beta test to the research and development community. About 40 universities and R&D laboratories in the US currently have the beta test code. The code release includes 200 pages of documentation and 15000 lines of C code and C shell scripts. The code has had three minor revisions since version 1.0. Version 2.0 is being prepared for release through COSMIC, and should be submitted October 1992. Version 2.0 includes extensions to the user interface and fixes all the bugs reported on the beta-test, but does not contain the decision graph routines. Version 3.0 is concurrently under development. This includes algorithms for learning decision graphs, and sophisticated any-time search algorithms for returning better quality trees and graphs. Version 3.0 is being released as beta test about November 1992.

Main use of the code to date has been in bench-marking, comparative studies, and comparative research on related algorithms, although groups in several different commercial and scientific areas currently have the code. Comparative studies done by several international research groups have found the code to be a good implementation,

somewhat slower than the original CART code, but more flexible, and easier to use. The new Bayesian extensions have also been found to give significant improvement over earlier tree algorithms, particularly in providing probability estimates, an important task for diagnosis and monitoring.

MODULES IN THE IND PACKAGE

The first task in using IND is to format your data into an appropriate text file and run it through the data conversion routine in IND. The routine `encsmpl` will produce an IND data description file for you, see Figure 4, and encode the data into IND's internal format. This data description file can then be modified to add defaults, utilities, constraints, etc., to configure IND for this data. An extract of a text file matching the description file in Figure 4 is given in Figure 5.

```

class : compensated_hypothyroid negative, primary_hypothyroid,secondary_hypothyroid.
age:      cont 0..100.
sex:      M,F.
on_thyroxine query_on_thyroxine on_antithyroid_medic sick pregnant thyroid_surgery l131_treatment : f,t.
TSH_measured: f,t.
TSH:      cont 0..600 (?).
TT4_measured: asfor TSH_measured.
TT4:      cont.
T4:       cont 0..3 (?).
FTI:      cont 0..400 (?).
referral_source: SVI,STMW,WEST,SVHC, SVHD,other (subset=full).

prior :    "-d8 -Anonym,1".
context :  TBG only if TBG_measured .

```

Annotations in the original figure:

- Arrow from `age` to `cont 0..100.`: this is the attribute to predict
- Arrows from `TSH` and `TT4` to `cont 0..600 (?)` and `cont.`: these attributes are identical types
- Arrow from `TT4` to `cont.`: missing values occur in this attribute
- Arrow from `T4` to `cont 0..3 (?)`: do subsetting on this attribute
- Arrow from `prior` to `"-d8 -Anonym,1".`: instructions to IND on default priors and constraints

Figure 4. The data description file input to IND.

```

negative 36 F f f f f f f t 0.22 t 191 0.98 194 other
negative 73 F f f f t f f f f ? t 119 0.92 129 SVI
compensated_hypothyroid 34 F f f f f t f f t 19 146 ? 125 other

```

Figure 5. Sample input data file matching description in Figure 4.

Once IND has the data encoded, IND can be operated at a number of different levels, depending on the requirements of the user. The simplest level is to use commands that have general prepared styles for tree generation. The command `mktree` shown in the top of Figure 6 uses prepared styles to drive the basic tree generation, pruning and classification routines. A sample run from `mktree` is given in Figure 7 at the end of the paper. This used the verbose option to automatically explain each component of IND and how it was configured. More experienced users of IND may like to make better use of the range of features. To do this, the low level routines can be called directly. All routines are controlled using the data description format of Figure 4 together with standard Unix style command options. Users may also wish to perform cross-validation to estimate error rates, or run experiments using a number of different tree styles to help in configuring IND for their problem. This can be done using the `ttest` utility shown at the top of Figure 6. This utility collates statistics required for you to analyze each run.

Some of the prepared styles available for the novice user of IND are as follows:

- `bayes, mml` : The simple.bayes style is useful when you know that most of the attributes supplied are relevant and that moderate accuracy is achievable. The mml style assumes poorer attribute quality. Both styles use Bayesian smoothing. These can also be modified with a look-ahead style.
- `cart` : A number of variations of basic CART are reimplemented in IND, although multivariate splits and surrogate splits are not implemented. Basic cart style using subsetting, twoing, cross validation cost complexity pruning and a simple stopping rule.
- `c4` : An early version of C4 is implemented with subsetting, pessimistic pruning and the gain ratio splitting rule.

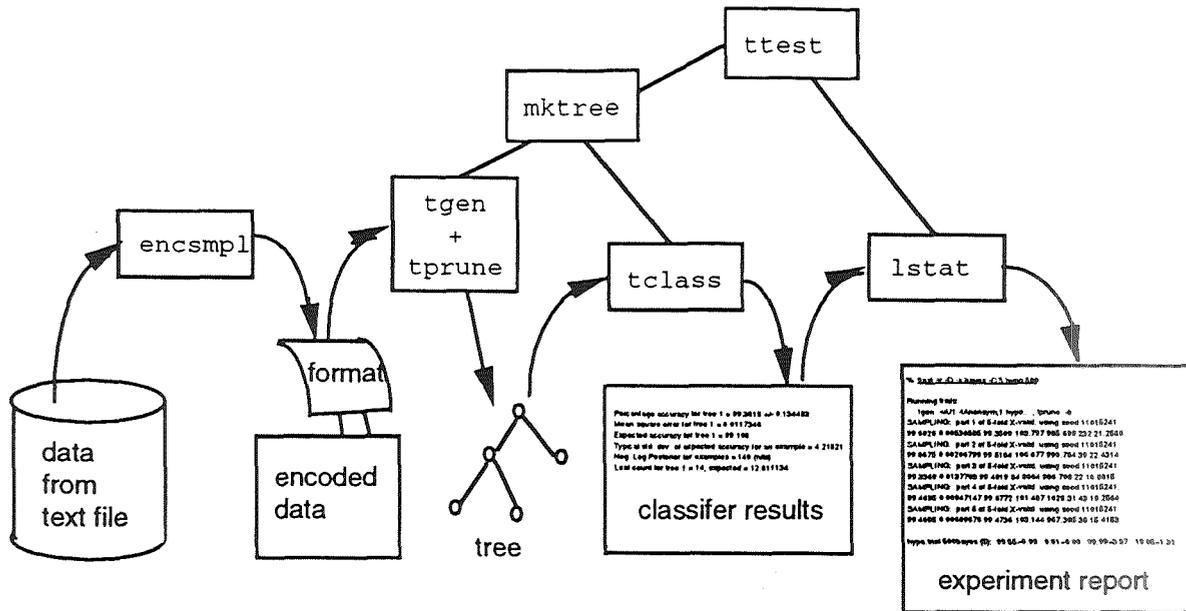


Figure 6. Overview of the modules in IND.

CLASS PROBABILITY TREE THEORY

In this section I briefly review the Bayesian theory of learning classification trees. This theoretical section should be skipped if your interests lie in applications of the algorithm. The section introduces the theory behind the unique Bayesian aspects of the IND package. More details of this theory are given in [5,6]. An excellent introduction to tree methods can be found in [2]. Theory behind the graph components of IND available in beta-test version 3.0 can be found in [7]. The methods discussed here are developed according to the algorithm design strategy presented in the earlier design section.

The basic tenet of Bayesian decision theory is that if we do not know something with reasonable certainty, then we should look at some reasonable and mutually exclusive alternatives and weigh them up, to help us make a "representative" decision. A reasonable alternative is one we currently have high subjective belief in. I will explain how this applies to trees, based on material in [6]. The formulation is sufficiently general so that it could just as well be applied to other classification models such as probabilistic rules, Bayesian networks, or one of many other knowledge representations from artificial intelligence, neural networks or statistics that have a probabilistic interpretation.

Class probability trees have a vector of class probabilities at their leaves, as shown in Figure 3c. They represent a conditional probability distribution of class value conditioned on other details about the case. A particular class probability tree can be represented by its discrete component T , the *tree structure* given by the shape of the tree and the tests at the leaves, and its continuous component S , the leaf class probabilities. This gives the conditional probability distribution $Pr(class|case, T, S)$, which is the *likelihood function* for a classified case (*class, case*) using the class probability tree specified by T and S .

Suppose we are given a training sample *Sample* consisting of classified cases *cases* and their classes *classes*, together with a new case, *new-case*, whose class, *new-class*, we wish to predict. If the goal is to minimize errors in prediction (other utility functions can be handled similarly), decision theory says we should choose the class *new-class* to maximize the posterior class probability $Pr(new-class|new-case, Sample)$. Using the tree model, this expression can be expanded using the laws of probability theory to obtain the posterior average of the class probabilities predicted for *new-class* from all possible class probability trees:

$$\begin{aligned}
 Pr(new-class|new-case, Sample) &= \sum_T \int_S Pr(new-class|new-case, T, S) Pr(T, S| Sample) dS \\
 &= \sum_T Pr(new-class|new-case, T, Sample) Pr(T| Sample)
 \end{aligned} \tag{1}$$

where the summations are over the space of all possible tree structures T , and

$$Pr(T | Sample) \text{ proportional-to } \int_S Pr(classes | cases, T, S) Pr(S | T) Pr(T) dS$$

$$Pr(new-class | new-case, T, Sample) \text{ proportional-to } \int_S Pr(new-class | new-case, T, S) Pr(S | T, Sample) dS$$

Formula (1) simply says to average the class predictions made for each tree. That is, since we aren't certain which tree is "true", we hedge our bets over reasonable trees. The *posterior* probability of the tree structure T , $Pr(T | Sample)$, is the weight used in the averaging process. The probabilities appearing in the formula above are calculated in log-space, to prevent underflow, and are sometimes referred to as "code-lengths" (because a negative log. probability is a code length by information theory).

The algorithm design strategy is based on designing a heuristic procedure to find a single tree or set of trees that can be used to approximate Formula (1). This is described by the following 4 steps.

Step 1. Develop priors over the structural and continuous components of the model, $Pr(S | T)$ and $Pr(T)$. The form of the prior should be flexible enough so that it can be changed from application to application. In the IND package, these priors can be tailored to your application, and advice is given in the manual. Alternatively, "bland" priors can be used if you don't wish to assume a particular prior.

Step 2. Given a training sample *Sample*, determine a suitably efficient way of computing or approximating the posterior of the structural component of the model. Then devise a heuristic search procedure for searching the space of structures to find structures with high posterior. In trees, a simple one-ply look-ahead procedure can be tried, which corresponds to the standard tree growing algorithm [2]. In IND, two-ply and three-ply versions of look-ahead are also available. These start with the trivial, empty tree. They then consider extending the tree by a single ply, by replacing an ungrown node with a test and leaves at its outcomes. A heuristic measure to evaluate the quality of a new growth can be determined from the posterior probabilities. Several different tests are tried and evaluated, and the best one is chosen for subsequent development.

Step 3. Given a training sample *Sample* and a structure T , determine a formula or approximation for the posterior expected values of the parameters S , $Pr(new-class | new-case, T, Sample)$, as required for Formula-(1).

Step 4. Devise a procedure for approximating the summation of Formula (1) by a small set of high posterior structures. There are three techniques for doing this:

Smoothing: The sum can be computed in closed form if it is restricted to the set of tree structures obtained by pruning a large tree structure in all possible ways. A linear time algorithm is given in [6]. This is called smoothing because it is equivalent to smoothing out the class probabilities at the leaf of a tree by averaging them the branch leading to the leaf. This is implemented in the "-b" option to IND's `tprune`.

Averaging: The sum can be approximated by searching for and storing many dominant terms, i.e. many high posterior trees structures. We can build multiple tree structures, and combine them together efficiently in an AND-OR representation called *option trees*. Growing option trees and then applying a similar summation process to smoothing is called *tree averaging*. This is implemented as a style in IND's `mktree`.

Multiple Models: The sum can be approximated by using importance sampling or Monte Carlo estimation. That is, a few tree structures are generated in approximate proportion with their posterior (this is done using the tree growing heuristic), and their class probability vectors uniformly averaged.

PERFORMANCE SUMMARY

Various experimental results from the use of IND version 1.0 are reported in [6]. Experimental results for the graph component of IND, available in beta-test version 3.0, can be got from results in [7] for earlier code from Jon Oliver. IND has been run on databases available from University of California at Irvine (FTP to `ics.uci.edu` and look in the directory `machine-learning-databases`). The results show that the new features of IND give more accurate class probability estimates for new examples, and often better predictions, though sometimes at the cost of increased computation, depending on the problem. The MML graph component of IND has previously been run by Oliver and colleagues on DNA structure data and produced results of interest to molecular biologists, see [7] and references therein for details. IND has recently been hooked up to the System Diagnostic Builder from GHG Corporation, which is used for building diagnosis systems at NASA's Johnson Space Center [8]

Acknowledgements

IND was based on an early suite of software developed at Basser Department of Computer Science at Sydney University by a lineage of students of Jason Catlett: David Harper, Murray Dean, David Muller and Chris Carter, and possibly some others. More recently Rich Caruana of CMU and Jon Oliver of Monash University worked on the package during summer internship at NASA-Ames Research Center. Also the users of the beta-release provided considerable feedback.

References

- [1] Berger, J. O. (1985). *Statistical Decision Theory and Bayesian Analysis*. Springer-Verlag, New York.
- [2] Quinlan, J. (1986). Induction of decision trees. *Machine Learning*, 1(1):81--106.
- [3] Quinlan, J. (1988). Simplifying decision trees. In Gaines, B. and Boose, J., editors, *Knowledge Acquisition for Knowledge-Based Systems*, 239--252. Academic Press, London.
- [4] Breiman, L., Friedman, J., Olshen, R., and Stone, C. (1984). *Classification and Regression Trees*. Wadsworth, Belmont.
- [5] Buntine, W. (1991). Classifiers: A Theoretical and Empirical Study. *International Joint Conference on Artificial Intelligence*, August, 1991, Sydney.
- [6] Buntine, W. (1992). Learning classification trees. *Statistics and Computing*, 2:63--73.
- [7] Oliver, J. (1992). Inferring decision graphs using the minimum message length principle. *Australian Artificial Intelligence Conference*, November, 1992, Australia.
- [8] Nieten, J. L. and Burke, R. (92). *System Diagnostic Builder*. Report from GHG Corporation at JSC.

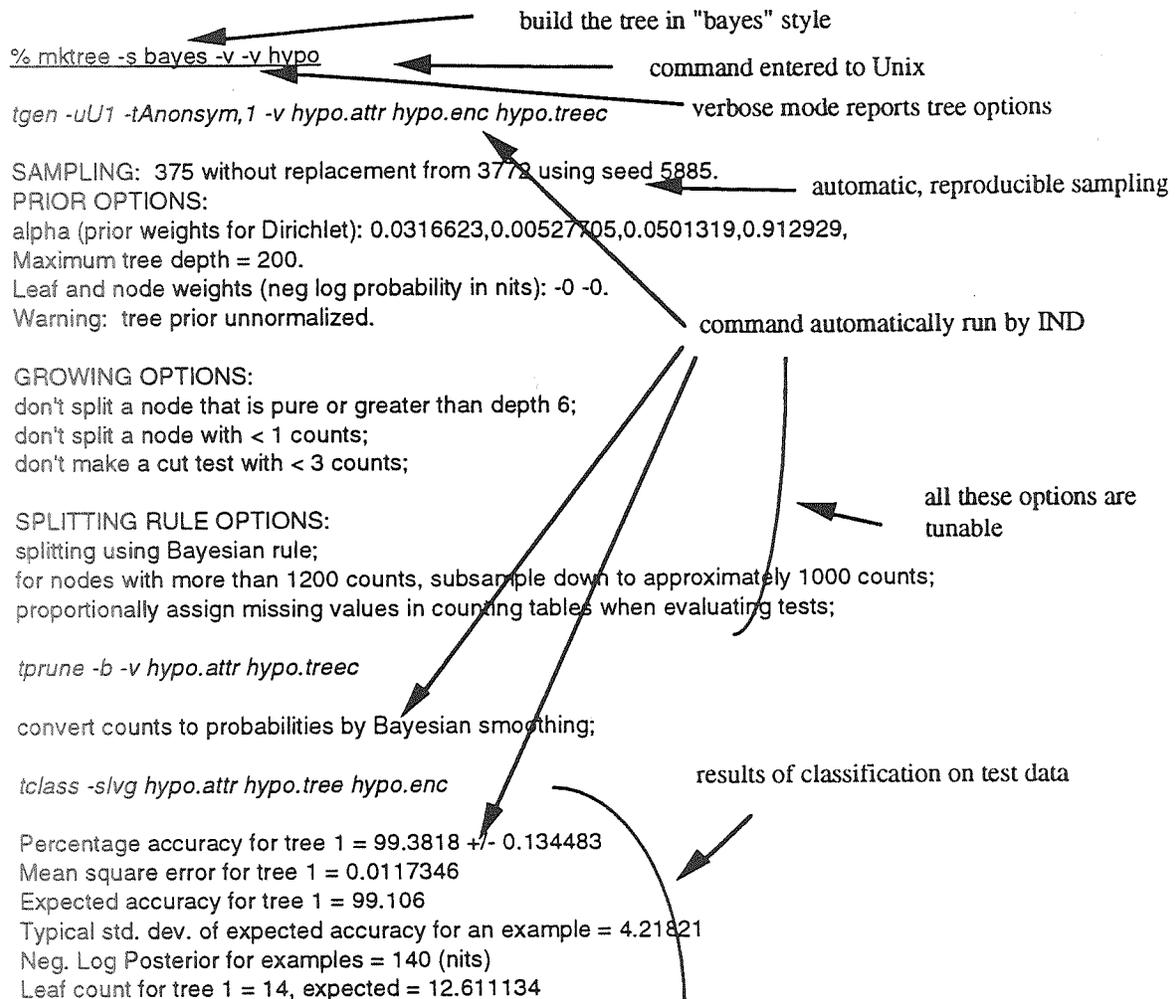


Figure 7. Building a tree using IND in verbose mode.

```
TSH < 6.05: 1.432e-05 0.0006207 2.909e-05 0.9993 negative
TSH >= 6.05:
| TSH_measured = f: 0.0001371 4.283e-06 0.0002784 0.9996 negative
| TSH_measured = t:
| | FTI < 64.5:
| | | T4_measured = f:
| | | | on_thyroxine = f:
| | | | | thyroid_surgery = f: 0.2523 9.126e-05 0.6358 0.1118 compensated_hypothyroid
| | | | | thyroid_surgery = t: 0.04943 0.0004665 0.1141 0.836 negative
| | | | on_thyroxine = t: 0.01588 0.0003925 0.03357 0.9502 negative
| | | | T4_measured = t:
| | | | | thyroid_surgery = f: 0.9637 1.226e-05 0.0007972 0.03548 primary_hypothyroid
| | | | | thyroid_surgery = t: 0.08208 0.0001835 0.01192 0.9058 negative
| | FTI >= 64.5:
| | | on_thyroxine = f:
| | | | TT4 < 150.5:
| | | | | thyroid_surgery = f: 0.1531 8.899e-05 0.7433 0.1035 compensated_hypothyroid
| | | | | thyroid_surgery = t: 0.00339 0.0001059 0.006886 0.9896 negative
| | | | TT4 >= 150.5: 0.04807 0.0001326 0.03691 0.9149 negative
| | | on_thyroxine = t: 0.0004708 1.471e-05 0.0009563 0.9986 negative
```

Figure 8. A print of the resultant tree showing class probabilities and decisions.

0 M12

**BIOTECHNOLOGY AND LIFE SCIENCES PART 2:
MEDICAL TECHNOLOGY**

299



AUTOMATIC DETECTION OF SEIZURES WITH APPLICATIONS

Dale E. Olsen, John C. Harris, Protagoras N. Cutchis, John A. Cristion
The Johns Hopkins University Applied Physics Laboratory
Johns Hopkins Road, Laurel, MD 20723-6099

Ronald P. Lesser, W. Robert S. Webber
The Johns Hopkins Hospital, Meyer Building
600 N. Wolfe St., Baltimore, MD 21287

ABSTRACT

There are an estimated two million people with epilepsy in the United States. Many of these people do not respond to anti-epileptic drug therapy. Two devices can be developed to assist in the treatment of epilepsy. The first is a microcomputer-based system designed to process massive amounts of electroencephalogram (EEG) data collected during long-term monitoring of patients for the purpose of diagnosing seizures, assessing the effectiveness of medical therapy or selecting patients for epilepsy surgery. Such a device would select and display important EEG events. Currently many such events are missed. A second device could be implanted and would detect seizures and initiate therapy. Both of these devices require a reliable seizure detection algorithm. A new algorithm is described. It is believed to represent an improvement over existing seizure detection algorithms because better signal features have been selected and better standardization methods have been used.

INTRODUCTION

Computerized analysis of electroencephalogram (EEG) data has been the goal of many investigators for years (refs 1-13, 15-28, 31). Such an analysis system could be used for diagnosing seizures, for assessing the effectiveness of medical therapy or for selecting patients for epilepsy surgery. A necessary component of accurate assessment of medical therapy, including safe and effective surgery, requires prolonged EEG monitoring to obtain recordings of seizure activity. One important reason for developing a system is to meet the needs of about one-third of the many patients with complex partial seizures who do not respond to anti-epileptic drug therapy. Computerized systems can be used to collect and store the data, but even with new mass-storage devices, it is not practical to store and review all of the data collected over days. Automatic spike and sharp wave detection and automatic seizure detection algorithms could be used to detect critical events and store just the relevant data for later examination by the neurologist and neural surgeons. This information, in conjunction with clinical observation, is used to localize epileptic form discharges. Many of the references [refs 4, 5, 7, 8, 13, 20] discuss the development of automatic spike and sharp wave detection. Reliable automatic seizure detection is perhaps more challenging and would be a key element of any system designed for analysis of EEG data.

Another use of an automatic seizure detection system derives from the research in the use of electrical stimulation to control epileptic seizures. Hammond, et al. [14], and Wilder, et al. [30], stimulated the vagus nerve, and Velasco, et al. [29], electrically stimulated the centromedian thalamic nucleus. Both groups significantly reduced clinical seizures. Zabara [31] produced, tested and patented an implanted device to control or prevent epileptic seizures. His device proved most effective for the patients who can feel a seizure coming on and then use a small magnet to initiate the stimulation of the vagus nerve. These reports, coupled with the possibility of drug therapy initiated using an implanted device, suggest that an implanted automatic seizure detection device, coupled with a therapeutic device, could be developed and used to significantly reduce the frequency, intensity or duration of seizures.

EXISTING DETECTION SYSTEMS

Five groups of investigators have pursued the goal of automatic seizure detection by implementing a system and publishing their results. Prior, et al. [refs 26, 27] developed a simple system aimed at recording the timing and frequency of seizure discharges. They were able to detect severe seizures by looking for large amplitude signals sustained over a period of time. Their system used a single channel of data recorded on paper running at 5 to 6 cm/hour.

Babb, et al. [1] developed electronic circuits for recording and detecting seizures. Their system used implanted electrodes to detect seizures which might otherwise go unreported. The system identified as seizure high-amplitude EEG signals of sufficient durations. They reported monitoring four patients and detecting 66 seizures, 20 of which were false detections. During this same period of time, the nurses detected 28 seizures, two of which were false detections. Therefore, the automated system detected 20 unreported seizures.

Ives, et al. [refs 15, 16] used implanted depth electrodes and a PDP-12 computer to remotely monitor patients. The system was adjusted or calibrated for individual patients using their seizure patterns. If the EEG signals fell within an amplitude window for a certain period of time, the algorithm identified the event as a seizure and recorded eight channels of EEG data. During two weeks of monitoring, thirteen seizures were detected and recorded by the computer. Only one seizure was reported by the nursing staff and only 3 were noted by the patient.

Gotman [9-13] designed a system to select, from largely uneventful EEG data, the sections which are likely to be of interest. Events, such as seizures, are recorded either because the program detects them or because the event button is pressed. His system was originally implemented in 1975 and has been in use since that time. It detects seizures using both surface and depth electrodes. His system is in use in many centers and has been updated [10]. A study of his latest algorithm used 5303 hours of data and showed that 24 percent of the 244 seizures recorded were missed by the automatic detection system. However, in 41 percent of the seizures, the patient alarm was not pressed, but the computer made the decision to record the event. Like the other systems, Gotman's system detects many seizures which would otherwise be missed. The system experiences about one false detection per hour of recording.

Gotman's method uses digitally filtered data broken up into 2-second epochs. He compares features of this data to what he calls "background." The background data is a 16-second long section of data ending 12 seconds prior to the epoch being analyzed. This comparison allows the algorithm to self-scale to account for differences in montage or other settings and is a form of what we will refer to as calibration or standardization. Detection occurs when 1) the average amplitude of half-waves in the epoch is at least three times that of the background, 2) their average duration corresponds to frequencies between 3 and 20 Hz, and 3) the coefficient of variation (ratio of the variance to the square of the mean) of the half-waves is below 0.36. In short the algorithm is based on amplitude, frequency, and the regularity of the half-wave duration. In 1990, Gotman [10] reported modifications to his algorithm including the requirement to look 8 seconds ahead to verify that the amplitude remains high. This modification reduced false detections.

More recently Murro, et al. [23] developed a system using concepts similar to those of Gotman but used spectral concepts for feature development. They used three features from each of two channels. These six features are reduced to four using principal component analysis, and then statistical discrimination is used to develop a detection rule. Their algorithm was developed for intracranial data only.

They address the issue of calibrating for patient differences partly in the way they define their features. They define reference power as the spectral power between 0.15 and 36 Hz averaged over four consecutive EEG epochs from 27 to 55.6 seconds prior to the event being evaluated. One of their features is relative power and is defined as the power between 0.15 and 36 Hz of the current epoch divided by the reference power. The second of their features is the dominant frequency. The third feature is rhythmicity, which is the ratio of the power associated with the dominant frequency to the relative power. Murro, et al.

use a segment of recent EEG data from the patient to influence the definition of their first and third feature. However, Murro, et al. also provide more in the way of custom tailoring their algorithm to the individual by building a separate decision rule for each patient. Their rule is based on normal data from the patient collected at different times and seizure data collected from other patients. These procedures allow for a more careful calibration for patient differences.

The algorithm of Murro, et al. [23] was tested using 8 patients and 43 seizures. Their system was evaluated using different detection thresholds. It detected all seizures allowing for a rate of 2.5 false detections per hour and detected 91% of the seizures with 1.5 false detections per hour.

ALTERNATIVE APPROACHES

The first seizure detection methods relied on amplitude and duration to identify seizures. Gotman's method tends to mimic the EEG readers. He characterizes the signals using features appearing to be motivated by those observed during seizure activity. His self-scaling and his detection rules are simple but effective.

The methods of Murro, et al. are similar to Gotman's but are more statistically sophisticated. Like Gotman, Murro, et al., used recent epochs to self-scale their features. They characterize the epochs of data using relative amplitude, dominant frequency and rhythmicity and then use statistical discriminate analysis to develop a detection rule.

There are three important common elements in these two algorithms and in our own:

- 1) standardization, calibration or self-scaling
- 2) feature selection
- 3) discrimination or a decision rule

We believe that the most critical element in discriminating seizure from other data is the definition and evaluation of features or feature selection. We also believe that other standardization or calibration methods will prove useful. Our algorithm has characteristics which it make fundamentally different from those of other investigators. More emphasis is placed on developing and evaluating a variety of features and on developing an improved approach to standardization.

THE ALGORITHM

Each channel or time series of EEG data is divided into non-overlapping 3-second epochs of data. The first goal is to compute a probability that the patient is having a seizure for each 3-second segment of data and for each of the different channels processed by the system. These probabilities will be combined with probabilities from adjacent epochs to determine if a patient is having a seizure.

To compute the probabilities, several steps are necessary. First the data must be standardized to account for patient, channel and hardware differences. The intent here is different from that of Gotman's self-scaling in that we are not simply evaluating the current epoch relative to background. The methods are very different and the results are different. However, both provide essential calibration necessary to account for patient differences. In our method, a standard deviation of the data is computed and updated for each channel. That is, values are recursively updated as new data enters the system. The current standard deviation is then used to scale the new data to make it consistent with that assumed by the algorithm. A standard deviation is computed and used for each channel.

The next step in developing probabilities of seizures is to characterize each epoch of data using many features and to evaluate these features to determine which can best help separate seizure and normal data. Several hundred features have been evaluated. These include usual time-series characteristics, features

which mimic EEG readers' methods and statistical characteristics. We find that all three types of features contribute to the discrimination and algorithms which use all three in combination are better able to detect seizures than those that do not.

Once a set of most-effective features has been identified, then any of the many discriminate techniques can be applied. Our base algorithm uses logistic regression. Neural networks and statistical discrimination techniques such as those described by Murro et al., [23] or Olsen et al., [24] could easily be applied. Logistic regression has the advantage that it provides a meaningful probability of seizure. In our algorithm, the user can adjust the probability levels to reduce false seizure detections or to increase sensitivity. The user may also adjust the standardization for the same purposes.

TRAINING AND TESTING THE ALGORITHM

The critical measure of an algorithm's performance is not how well it detects seizures, but how well it detects seizures without producing false detections. It is a simple matter to demonstrate an algorithm correctly identifies seizures on a limited span of data. Our algorithm has been applied to several spans of seizure data from different patients and all seizures were detected with no false detections. The difficulty comes when extending an algorithm's use to days of data. Thus, while our algorithm performed admirably on many data sets, the only true test is that which comes from extensive testing on large representative data sets or from continuous monitoring of patients in a real-time environment.

Five-minute segments of data were collected from each of two patients, both day and night, over a five-day period. The segments were sampled at least once per hour with an effort being made to record a variety of activity and sleep stages. Nine seizures (some lasting over 10 minutes) were recorded, five from one patient and four from the other. Seven and one-half hours of sampled data were saved. An initial algorithm was developed for two channels from one patient and was successfully applied to the second. All seizures in both patients were detected and there were no false detections.

Our algorithm was developed using all 35,000 epochs of standardized data from both patients. Seizure onset for one patient was on the left temporal region and for the second was on the right. As a result, data from FP1-F7, F7-T3, T3-T5, and T5-O1 from the first patient and data from FP2-F8, F8-T4, T4-T6, and T6-O2 from the second were used as example data to build the algorithm. When the algorithm was applied to the sampled data, *all seizures were detected and there were no false detections.*

Our algorithm used logistic regression to identify the most useful features and to form a decision rule. Unlike neural networks, logistic regression selects only about 5-15 constants and therefore cannot over-train or memorize even with much smaller data sets. As a result, the testing on 35,000 epochs of data described above provides significant evidence of the discrimination capability of the algorithm. It remains to test the algorithm on other large data sets and in a real-time environment.

APPLICATIONS OF AUTOMATIC SEIZURE DETECTION

In the United States alone, there are over two million people with epilepsy. Each year 100,000 new cases are diagnosed. The most common type of seizures are complex partial seizures and of these, 35% of the patients fail to respond to anti-epileptic drug therapy [14]. Since for these patients, long-term monitoring is necessary for localization of epileptic activity (for possible local resection), the need is clear for a system which will process the massive EEG data collected and decide which data to save. Such a system would also be used to evaluate various drug therapies. Automatic spike and seizure detection capability is an essential part of developing a microcomputer based monitoring system.

As noted in the introduction, the first application of the seizure detection algorithm will be to monitor patients in long-term epilepsy monitoring units. Most of these data will be from surface electrodes.

A more ambitious application, in some respects, is to automatically detect seizures and reduce their intensity or duration. Since the data for this application would be from implanted electrodes, the seizure detection algorithm would not need to deal with muscle and other artifacts or the effects of the skull on the electrical signals. The algorithm could be adapted to the patient. Based on experience with more difficult surface detection algorithms, we believe that a highly reliable detection algorithm can be developed.

SEIZURE DETECTION IMPLANT HARDWARE

It is somewhat difficult at this stage of the program to predict accurately what the configuration of an implanted monitoring device will be. This is due to the fact that the algorithm would not be finalized until extensive on-line testing is complete. However, the development of such a device is not beyond the "state of the art" and certain aspects of the implant can be predicted, which are discussed here. The physical configuration will be very much like a modern-day pacemaker. There will be a titanium case about 1.5" x 1.5" x 0.3" thick that will be implanted subcutaneously just below the clavicle. This case will hold the electronics as well as a lithium battery which should allow an implant life of 5-6 years. There will be a lead containing at least 4 conductors connecting to the case. This lead will be tunneled subcutaneously over the clavicle, up the neck and under the scalp to the point of entry into the skull. This is the same technique which is currently used for hydrocephalus shunt systems. Once the leads enter the skull, they will travel under the skull to the surface locations on the brain which were selected by the surgical team at implant.

The circuitry for this device will probably use subthreshold CMOS analog circuitry rather than digital CMOS circuitry. The analog circuitry uses less power for the large number of computations which will be required by this device. The analog approach also eliminates the analog-to-digital (A/D) converters which would be required in the digital approach.

CONCLUSION

The development of new medical procedures has enabled neurologists and neural surgeons to provide safe and effective treatment of otherwise intractable epilepsy. At the same time, the development of computer technology has made it possible to collect and process more information than ever before. Key elements to improving the general care of patients with intractable epilepsy are the algorithms which will enable the computers to efficiently reduce the data collected to information useful in planning therapy or possibly to initiate treatment in a timely fashion. Systems for the long-term EEG monitoring have been developed at the Johns Hopkins Hospital and elsewhere. Some of these systems automatically save spike and seizure data for later analysis by neurologists. Commercial development of such a system is possible.

REFERENCES

- ¹ Babb, T. L., Mariani, E., Crandall, P. H., "An Electronic Circuit for Detection of EEG Seizures Recorded with Implanted Electrodes," *Electroencephalograph and Clinical Neurophysiology*, 37, 305-308 (1974).
- ² Barlow, J. S., "Methods of Analysis of Nonstationary EEGs, with Emphasis on Segmentation Techniques; A Comparative Review," *J. Clin. Neurophysiol.* 2, 267-304 (1985).
- ³ Bender, R., Schultz, B., Schultz, A., Pichlmayr, I., "Identification of EEG Patterns Occurring in Anesthesia by Means of Autoregressive Parameters," *Biomedizinische Technik*. 36, 236-240 (1991).

- ⁴ Frost, J. D., Jr., "Automatic Recognition and Characterization of Epileptiform Discharge in the Human EEG," *J. Clin. Neurophysiol.* 2, 231-249 (1985).
- ⁵ Frost, J. D., Jr., "Microprocessor-based EEG Spike Detection and Quantification," *Int. J. Bio-Med. Comput.* 10, 375-373 (1979).
- ⁶ Gevins, A. S., Yeager, C. L., Diamond, S. L., Spire, S.P., and Zeitlin, G.M., "Automated Analysis of the Electrical Activity of the Human Brain (EEG): A Progress Report," *Proc. IEEE* 63, 1382-1399 (1975).
- ⁷ Gevins, A. S., Blackburn, J., and Dedon, M., "Very Accurate Computer Recognition of Three-per-Second Generalized Spike-and-Wave Discharge," in *Advances in Epileptology; the Xth Epilepsy International Symposium*, Wada, J. A., and Dentry, S. K. (eds.), Raven Press, N.Y., 121-128 (1980).
- ⁸ Glover, J. R., Ktonas, P. Y., Raghavan, N., Urunvela, J. M., Velamuri, S. S., Reilly, E., "A Multichannel Signal Processor for Detection of Epileptogenic Sharp Transients in the EEG," *IEEE Transactions on Biomedical Engineering*, BME-33, No. 12, 1121-1128 (1986).
- ⁹ Gotman, J., and Gloor, P., "Automatic Recognition and Quantification of Epileptic Activity in the Human Scalp (EEG)," *Electroencephalogr. Clin. Neurophysiol.* 48, 551-557 (1980).
- ¹⁰ Gotman, J., "Automatic Seizure Detection: Improvements and Evaluations," *Electroencephalograph and Clinical Neurophysiology*, 76 317-324 (1990).
- ¹¹ Gotman, J., "Computer Analysis During Intensive Monitoring of Epileptic Patients," *Advances in Neurology*. 46, in *Intensive Neurodiagnostic Monitoring*, R. J. Gumnit (ed.), Raven Press, New York (1986).
- ¹² Gotman, J., "Practical Use of Computer-Assisted EEG Interpretation in Epilepsy," *J. Clin. Neurophysiol.* 2, 251-265 (1985).
- ¹³ Gotman, J., "Quantitative Measures of Epileptic Spike Morphology in the Human EEG," *Electroencephalogr. Clin. Neurophysiol.* 48, 551-557 (1980).
- ¹⁴ Hammond, E. J., Uthman, B. M., Reid, S. A., Wilder, B. S., Ramsay, R. E., "Vagus Nerve Stimulation in Humans: Neurophysiol Studies and Electrophysiological Monitoring," *Epilepsia* 31, 51-59 (1990).
- ¹⁵ Ives, J. R., Thompson, C. S., Gloor, P., Olivier, A., Woods, S. F., "The On-line Computer Detection and Recording of Spontaneous Temporal Lobe Epileptic Seizures from Patients with Implanted Depth Electrodes Via a Radio Telemetry Link." *Electroencephalography and Clinical Neurophysiology*, 37, p. 205 (1974).
- ¹⁶ Ives, J. R., Thompson, C. J., Gloor, P., "Seizure Monitoring: A New Tool in Electroencephalography," *Electroencephalography and Clinical Neurophysiology*, 41, 422-427 (1976).
- ¹⁷ Jansen, B. H., "Analysis of Biomedical Signals by Means of Linear Modeling," *CRC Crit. Rev. Biomed Eng.* 12 (4), 343-392 (1985).
- ¹⁸ Jansen, B. H., Hesman, A., and Lenten, R., "Piecewise Analysis of EEGs Using AR-Modeling and Clustering," *Comput. Biomed. Res.* 14, 168-178 (1981).
- ¹⁹ Jansen, B. H., Bourne, S. R., and Ward, S. W., "Autoregressive Estimation of Short Segment Spectra for Computerized EEG Analysis," *IEEE Trans. Biomed. Eng.* 26, 630-637 (1981).

- ²⁰ Ktonas, P. Y., "Automatic Spike and Sharp Wave (SSW) Detection." in *Methods of Analysis of Brain Electrical and Magnetic Signals EEG Handbook*, A. S. Gevins and A. Remond (eds.) Elsevier Science Publishers, 212-241 (1987).
- ²¹ Lopes Da Silva, H. A., Disk, A., and Smith, H., "Detection of Nonstationarities in EEGs Using the Autoregressive Model--An Application of EEGs of Epileptics." in *CLEAN: Computerized EEG Analysis*, Dolce, G., and Kunkel, H. (eds.), Gustav Fischer Verlag, Stuttgart, 180-199 (1975).
- ²² Morf, M., Vierira, A., Lee, D. T., and Kailath, T., "Recursive Multichannel Maximum Entropy Spectral Estimation," *IEEE Trans. Geosci. Electron.* 16, 85-94 (1978).
- ²³ Murro, A. M., King, D. W., Smith, J. R., Gallagher, B. B., Flanigin, H. F., Meador, K., "Computerized Seizure Detection of Complex Partial Seizures," *Electroencephalography and Clinical Neurophysiology*, 79, 330-333 (1991).
- ²⁴ Olsen, D. E., Criston, J. A., Spaur, C. W., "Automatic Detection of Epileptic Seizures Using Electroencephalographic Signals," *The Johns Hopkins APL Technical Digest*, 12, 182-191 (1991).
- ²⁵ Penczek, O., Grochulski, W., Greya, J., and Kowalczyk, M., "The Use of Multi-Channel Kalman Filter Algorithm in Structural Analysis of the Epileptic EEG," *Int. J. Bio-Med. Comput.* 20, 135-151 (1987).
- ²⁶ Prior, P. F., Virden, R. S. M., and Maynard, D. E., "An EEG Device for Monitoring Seizure Discharges," *Epilepsia* 14, Elsevier Science Publishers, 367-372 (1973).
- ²⁷ Prior, P. F., Maynard, D. E., "Recording Epileptic Seizures," *Proceedings of the Seventh International Symposium of Epilepsy, Berlin (West)*, Thieme Edition, Publishing Sciences Group, Inc., 325-377 (1975).
- ²⁸ Sato, K., Ono, K., Chiba, G., and Fukata, K., "On Some Methods for EEG Pattern Discrimination," *Int. J. Neurosci*, 7, 201-206 (1987).
- ²⁹ Velasco, F., Velasco, M., Ogarrio, C., Fanghanel, G., "Electrical Stimulation of the Centromediam Thalamic Nucleus in the Treatment of Convulsive Seizures; A Preliminary Report," *Epilepsia*, 28 (4), 421-430 (1987).
- ³⁰ Wilder, B. J., Uthman, B. M., Hammond, E. J., "Vagal Stimulation for Control of Complex Partial Seizures in Medically Refractory Epileptic Patients," *Pacing and Clinical Electrophysiology*, Future Publishing Co. 108-115 (1991).
- ³¹ Zabara, J., Neurocybernetic Prosthesis, Patient Number 5,025,807, June 25, 1991.

532-52
150502
p. 9

A FIBER OPTIC PROBE FOR THE DETECTION OF CATARACTS -

N 9 3 - 2 5 5 9 3

Rafat R. Ansari
NASA Lewis Research Center/CWRU, Mail Stop 105-1,
21000 Brookpark Road, Cleveland, Ohio 44135

Harbans S. Dhadwal
Department of Electrical Engineering
State University of New York at Stony Brook
Stony Brook, New York 11794

ABSTRACT

A compact fiber optic probe developed for on-orbit science experiments has been applied to detect the onset of cataracts, a capability that could eliminate physicians' guesswork and result in new drugs to "dissolve" or slow down the cataract formation before surgery is necessary. The probe is based upon dynamic light scattering (DLS) principles. It has no moving parts, no apertures, and requires no optical alignment. It is flexible and easy to use. Results are presented for excised but intact human eye lenses. In a clinical setting, the device can be easily incorporated into a slit-lamp apparatus (ophthalmoscope) for complete eye diagnostics. In this set up the integrated fiber optic probe, the size of a pencil, delivers a low power cone of laser light into the eye of a patient and guides the light which is back scattered by the protein molecules of the lens through a receiving optical fiber to a photo detector. The non-invasive DLS measurements provide rapid determination of protein crystalline size and its size distribution in the eye lens.

INTRODUCTION

A normal eye lens is a jelly-bean-size transparent tissue. A cataract is formed when this lens becomes cloudy. This cloudiness or opacification hinders light transmission through the lens and the ability to focus a sharp image on the retina at the back of the eye. The common symptoms experienced with cataracts include blurred or double vision, sensitivity to light and glare, less vivid perception of color, and frequent eyeglass prescription changes. If a large portion of the lens becomes cloudy, sight can be partially or completely lost until the cataract is removed¹. At present there is no medical treatment which will prevent cataracts or reverse them once they develop. The only treatment is the surgical removal of the clouded lens and its replacement with an intraocular lens implant. Today, an estimated 1.4 million cataract surgeries are performed each year in the United States. In normal patients, cataracts develop gradually over a period of many years. A cataract is caused by a change in the chemical composition of the lens. These changes can be attributed to aging (senile cataract), eye injuries (traumatic cataract), certain diseases and conditions of the eye and body (secondary cataract) e.g. high blood-sugar levels in diabetic patients, and hereditary or birth defects (congenital cataract). The senile (age-related) cataracts are the most common type of cataracts. They can occur as early as age 40.

An adult eye lens is primarily comprised of about 65 wt.% water and 35 wt.% proteins, the highest of any tissue in the body. The protein molecules or crystallines in the lens are subdivided into α , β , and γ crystallines. The α , β , and γ crystallines have a molecular weight of $\sim 10^6$ daltons, $\sim 10^5$ daltons, and $\sim 2 \times 10^4$ daltons respectively. Since α crystallines are the largest molecules, they are the strong scatterers of light in a light scattering experiment. When these protein molecules are agglomerated, they give rise to lens opacities. This has been confirmed by a variety of complementary techniques ranging from biochemical², light scattering³⁻⁵, QELS⁶⁻¹¹ and electron microscopy¹². Current clinical apparatus, which include visual inspection through a slit lamp microscope, and analysis of a photographic plate, lack the sensitivity and accuracy to detect small cellular and biochemical changes¹³. Biochemical studies of lens extract have demonstrated that aging in the normal mammalian lens is accompanied by conversion of the α -crystalline into higher molecular weight species¹⁴⁻¹⁶. In earlier stages of cataract formation the patients' cataracts are difficult for a physician to diagnose because of the lack of non-invasive experimental techniques. In general, progressively deteriorating vision conditions and in particular frustrating night vision conditions are unfortunately the only indications of a semi-developed cataract. This poses serious emotional situations for some patients and

forces guesswork on the physician's part for making well informed decisions concerning performance of surgical procedures.

In this paper we report the application of a fiber optic probe for the early detection of cataracts in excised, but intact human eye lenses, and its incorporation in a clinical set up. Earlier investigations have been reported in our recent publications^{17,18}. We believe these new technological advances will be useful in the investigation of cataracts during their early stages of formation. This could result in eliminating a physician's guess-work, reducing a patient's emotional trauma level, and this may allow pharmaceutical companies or dietitians to formulate new drugs or food supplies to "dissolve" the cataract, hence reducing the risks of costly and unwanted surgical procedures. Looking ahead to the year 2000, it is anticipated that the evolution of new drugs will slow the development of lens opacities that are leading cause of blindness worldwide¹⁹.

Dynamic Light Scattering (DLS)

DLS is an established laboratory technique which provides non-invasive measurements of particle size and size distributions, molecular weight, and particle-particle interactions for particles suspended in dilute solutions. The range of application ranges from 3 nm to 3 μm . Several books are available on this subject²⁰ (see reference 20 and the references contained therein). The data analysis techniques in DLS are reported in a review article²¹. Several other names are frequently used for DLS. These are Quasi-Elastic Light Scattering (QELS), Photon Correlation Spectroscopy (PCS), and Intensity Fluctuation Spectroscopy (IFS). In these experiments the laser light is focussed into a small spot inside the sample. The scattering volume, defined by the intersection of the incident and detection geometries, normally contains submicron particles suspended in a fluid medium. The intensity of the scattered light fluctuates due to the thermal movement (Brownian motion) of the particles. The intensity fluctuations in the scattered light are detected by a photodetector. This detected signal is processed via a digital correlator to yield an autocorrelation function. For dilute dispersions of spherical particles the slope of the autocorrelation function provides a quick and accurate determination of the particle's translation diffusion coefficient, which can be related to its size via a Stokes-Einstein equation, provided the viscosity of the suspending fluid, its temperature, and its refractive index are known. For concentrated suspensions and for dispersions containing more than one scattering species in the suspending medium, however, the data analysis and interpretation of the autocorrelation function becomes more difficult. We have discussed some of these data analysis techniques in our prior publications^{17,18,22,23}. The in vivo uses of conventional DLS technique as applied to the anterior segment of the eye have been reviewed elsewhere¹¹.

Until recently conventional DLS instruments have relied on bulky laser sources, bulk optics, and high voltage detection devices e.g. photomultiplier tubes (PMT). The incorporation of new advances in solid state technology, and the development of compact DLS spectrometers (fiber optic probes), have paved the way for a next generation laser light scattering instrument (LLSI). More recently, a compact, rugged, and modular LLSI has been conceived for microgravity experiments on board the space shuttle orbiter and possibly space station Freedom to study a range of phenomena²⁴. These phenomena include, among others, nucleation in crystal growth, aggregation, polymer induced flocculation, gelation, critical phenomena, and spinodal decomposition.

The most promising clinical technique that has become available is based on quasielastic light scattering (QELS). Early work establishing the utility of QELS for detection of molecular changes in the lens were demonstrated by Tanaka and Benedek²⁵. A recent clinical QELS study²⁶ of diabetic and non-diabetic patients has attempted to correlate the QELS results with visual inspection. Early researchers were concerned with solving the mystery concerning the transparency of a normal adult lens, given the high concentration of proteins in an aqueous solution²⁷⁻³⁰. A reliable, quantitative technique, causing the least trauma to the patient, has been a long sought goal for the study of cataractogenesis and other ocular disorders. Although the technique of QELS or DLS was first applied to study cataractogenesis by Tanaka and Benedek many years ago²⁵, its commercial scope has remained limited because of elaborate instrumentation, bulk optics and associated optical alignment problems, statistical errors in data analysis,

multiple scattering problems associated with mild and severe cataracts, and the polydisperse nature of the cataract itself.

EXPERIMENTAL PROCEDURE

Fiber optic probe

Our fiber optic probe employs two monomode optical fibers. The fibers have a core radius of about 4 μm . One fiber is used to transmit a Gaussian laser beam to the scattering volume. The second fiber positioned at a backscatter angle of 155° acts as a scattered light receiver. In concentrated dispersions, we have found that the backscatter regime allows recovery of particle sizes well beyond the cutoff point for the conventional state of the art laser light scattering systems, without requiring any additional corrections for multiple light scattering²⁸. This feature is useful in the study of cataractogenesis because cataractous conditions may give rise to multiple scattering effects. Depending upon the severity of the cataract, multiple light scattering effects will introduce high frequency components in the light scattering spectrum. This may cause a loss of resolution and a subsequent broadening of the particle size distribution. The fiber optic probe used in this study alleviates these problems. The detailed design considerations of the fiber optic probe and its range of application to investigate concentrated dispersions have been published elsewhere²⁸. Compared to conventional commercial state-of-the-art DLS spectrometers, our probe is 1-2 orders of magnitude smaller in physical size, and is inexpensive to fabricate.

Experimental set up

The DLS system employed in this work is schematically represented in Figure 1. Our integrated probe is comprised of two optical fibers, which are positioned in close proximity to each other and mounted into a single stainless steel ferrule. A monomode optical fiber pigtailed to a semiconductor laser, guides a Gaussian laser beam to a point inside the eye lens. The optical fiber is ruggedized by threading it through a teflon tubing and an outer plastic monocoil tubing. A bare portion of this monomode optical fiber is epoxied into a precision machined hole. A second optical fiber, is positioned into another machined hole in close proximity to the transmitting optical fiber, and is used for coherent detection of the light scattered from the eye lens in the backward direction. The receiving optical fiber is threaded through the same teflon sleeve and monocoil tubing up to the point where the transmitting and receiving fibers are separated. The receiving fiber terminates in a connector which is mated directly with a miniature photomultiplier tube.

Clinical Set up

Figure 2 shows a schematic of the clinical apparatus which can be used for measuring the intensity autocorrelation from patients. The integrated probe shown in Figure 1 is mounted on a slit-lamp apparatus: an instrument of choice for the ophthalmologists. The position of the probe can be adjusted using the standard joystick; control lever for horizontal and vertical movement, available on the slit-lamp apparatus. This arrangement provides precise positioning and location of the scattering volume in any substantially transparent region of the anterior segment of the eye. The probe tip is positioned so that a point inside the patients' eye lens is illuminated with an expanding Gaussian laser beam, having a diameter of 4 μm . The second optical fiber collects the scattered light at a fixed scattering angle and is connected to a miniature photomultiplier tube, followed by a compact data acquisition system. The photon pulse train after suitable amplification and discrimination is correlated using a digital correlator in a lap-top computer.

The application of the fiber optic probe for the measurement of cataractogenesis is procedurally similar to techniques familiar to ophthalmologists, most notably applanation tonometry and ultrasonography. The procedure can be easily done at a slit lamp under the installation of topical (drop) anesthesia. The eye movement is negated by having the patient direct vision in the contralateral eye on the fixation light. Findings from other clinical systems³¹⁻³⁴, introducing laser beams into patients' eyes, have concluded that a 20 minute exposure is the upper limit before patient fatigue becomes a limiting factor. Our measurements should take less than 2 minutes per autocorrelation function i.e. one test per patient. The maximal retinal

irradiance is a function of wavelength, incident power, numerical aperture of the cone of laser light and exposure time. In our fiber optic system we expect retinal irradiance to be less than 0.05 mW/mm^2 , which is three orders of magnitude below the damage threshold of 2 W/cm^2 for a 10 second exposure³⁵⁻³⁸.

RESULTS

Figure 3 summarizes the results of this investigation on excised, but, intact human eye lenses. The five pairs of cadaver human eye lenses employed in this study belonged to 18, 43, 55, 65, and 73 year old patients. The crystalline size increase as a function of patients' age is consistent with the development of senile cataract in the middle age to golden age patients. Upon visual examination, by a professional ophthalmologist (Dr. M.A. Dellavecchia of the Bryn Mawr Eye Clinic in Philadelphia), the lenses of younger patients were found to be transparent while the older lenses seemed to have a yellowish tint in them, consistent with senile cataractous changes. Clinically, these older lenses can be classified as having a mild to moderate cataract. None of the lenses were completely opaque.

CONCLUSION

In this paper we have shown the application of a fiber optic probe to non-invasively detect the onset of cataracts in human eye lenses, and its incorporation into a slit-lamp apparatus for complete eye diagnostics of the anterior chamber of the eye. The probe has a unique design which does not require any lenses, has no moving parts, does not need alignment, and is insensitive to vibrations and RF interference. In a clinical environment this new capability, when used in conjunction with regular eye examination on an yearly basis, will let the physician detect an incipient cataract before it forms.

REFERENCES

1. see "Cataract: Clouding The Lens of Sight" 1984 American Academy of Ophthalmology publication, printed 11/1990.
2. J. A. Jedziniak, J. H. Kinoshita, E. M. Yates, L. O. Hocker, and G. B. Benedek, "On the presence and mechanism of formation of heavy molecular weight aggregates in human normal and cataractous lenses," *Exp. Eye. Res.*, **15**, 185-189 (1973)
3. J. A. Jedziniak, D. F. Nicoli, H. Baram, and G. B. Benedek, "Quantitative verification of the existence of high molecular weight protein aggregates in the intact normal human lens by light scattering spectroscopy," *Investigative ophthalmology*, **17**(1), 51-57 (1978)
4. A. P. Bruckner, "Picosecond light scattering measurements of cataract microstructure," *App. Opt.*, **17**(19), 3177-3183 (1978)
5. F. A. Bettelheim and M. Paunovi, "Light scattering of normal human lens I: Application of random density and orientation fluctuation theory," *Biophys. J.*, **26**, 85-100 (1979)
6. C. Andries, W. Guadens and J. Clauwert, "Photon and fluorescence correlation spectroscopy and light scattering of eye lens proteins," *Biophys. J.*, **43**, 345-354 (1983)
7. J. N. Weiss, S. E. Bursell, R. E. Gleason and B. H. Eichold, "Photon correlation spectroscopy of in vivo human cornea," *Cornea*, **5**(1), 19-24 (1986)
8. C. Andries and J. Clauwert, "Photon correlation spectroscopy and light scattering of eye lens proteins at high concentrations," *Biophys. J.*, **47**, 591-605 (1985)
9. P. C. Magnante, L. T. Chylack and G. B. Benedek, "In vivo measurements on human lens using quasielastic light scattering," *SPIE Proceeding*, **605**, 94-97, (1986)
10. M. Delaye, J. I. Clark, and G. B. Benedek, "Identification of the scattering elements responsible for lens opacification in cold cataracts," *Biophys. J.*, **37**, 647-656 (1982)
11. S. E. Bursell, P. C. Magnante and L. T. Chylack Jr., "In vivo uses of quasielastic light scattering spectroscopy as a molecular probe in the anterior segment of the eye," in *Noninvasive Diagnostic Techniques in Ophthalmology*, editor Barry R. Masters, Springer-Verlag, New York (1990)
12. F. A. Bettelheim, E. L. Siew, S. Shyne, P. Farnsworth and P. Burke, "A comparative study of human lens by light scattering and scanning electron microscopy," *Exp. Eye. Res.*, **32**, 125-129 (1981)
13. A. Spector, S. Li and J. Sigelman, "Age-dependent changes in the molecular size of human lens proteins and their relationship to light scatter," *Investigative ophthalmology reports*, **13**(10), 795-798 (1974)
14. H. F. Honders and H. Bloemendal, "Aging of lens proteins" in *Molecular and Cellular Biology of the Eye Lens*, editor H. Bloemendal, 229-236, Wiley Interscience, New York (1981)
15. L. J. Takemoto and P. Azari, "Isolation and characterization of covalently linked, high molecular weight proteins from human cataractous lens," *Exp. Eye. Res.*, **24**, 63-70 (1977)
16. M. J. McFall-Ngai, L. L. Ding, L. J. Takemoto and J. Horwitz, "Spatial and temporal mapping of age-related changes in human lens crystallins," *Exp. Eye. Res.*, **41**, 745-758 (1985)
17. R. R. Ansari, H. S. Dhadwal, M. C. W. Campbell and M. A. Dellavecchia, "A fiber optic sensor for ophthalmic refractive diagnosis," *SPIE vol. 1648, Fiber Optic Medical and Fluorescent Sensors and Applications*, 83-105, Optical Engineering Press, Washington (1991).
18. H.S. Dhadwal, R.R. Ansari, and M.A. Dellavecchia, "A Coherent Fiber Optic Sensor for the Early Detection of Cataractogenesis in a Human Eye Lens", Accepted for publication in *J. Opt. Engineering: Biomedical Optics*, February 1993.
19. M.E. Long, "The Sense of Sight", in *National Geographic*, **182** No. 5., Nov. 1992.
20. B. Chu, editor, "Selected Papers on Quasielastic Light Scattering by Macromolecular, Supramolecular, and Fluid Systems", *SPIE Milestone Series, Volume MS 12*, SPIE Optical Engineering Press, Bellingham, Washington, 1990.
21. R.S. Stock, W.H. Ray, "Interpretation of Photon Correlation Data: A comparison of Analysis Methods", *J. Polym. Sci.: Polymer Physics Edition*, **23**, 1393-1447 (1985).
22. H.S. Dhadwal, R.R. Ansari and W.V. Meyer, "A Fiber Optic Probe for Particle Sizing in Concentrated Suspensions", *Rev. Sci. Instruments*, **60**, 12, 1991.
23. H.S. Dhadwal, R.R. Ansari, "Multiple Fiber Optic Probe for Several Sensing Applications", *Fiber Optic and Laser Sensors IX, SPIE Proceedings, 1584*, 262-272, 1991, Boston, MA.

24. W.V. Meyer, R.R. Ansari, "A Preview of a Microgravity Laser Light Scattering Instrument", AIAA 91-0779, 29th Aerospace Sciences Meeting, January 7-10, 1991, Reno, NV.
25. T. Tanaka and G. B. Benedek, "Observation of protein diffusivity in tact human and bovine lenses with application to cataract," *Investigative ophthalmology*, **14**(6), 449-456 (1976)
26. S. Bursell, R. S. Baker, J. N. Weiss, J. H. Haughton and L. I. Rand, "Clinical photon correlation spectroscopy evaluation of human diabetic lenses," *Exp. Eye. Res.*, **49**, 241-258 (1989)
27. S. Trokel, "The physical basis for transparency of crystalline lenses," *Investigative ophthalmology*, **1**, 493-501 (1962)
28. D. M. Maurice, *J. Physiol. (London)*, **736**, 263-268 (1957)
29. G. B. Benedek, "Theory of transparency of the eye," *App. Opt.*, **10**(3), 459-473 (1971)
30. M. Delaye and A. Tardieu, "Short-range order of crystallin proteins accounts for eye lens transparency," *Nature*, **302**, 415-417 (1983)
31. C.E. Riva, G.T. Fake, B. Eberly, and V. Benary, "Bidirectional LDV System for Absolute Measurement of Blood Speed in Retinal Vessels", *App. Opt.*, **18**, 2301-2306 (1979)
32. B.L. Petrig and C.E. Riva, "Retinal Laser Doppler Velocimetry: toward its computer-assisted clinical use", *App. Opt.*, **27**, 1126-1134 (1988)
33. B.L. Petrig and C.E. Riva, "Near-IR Retinal Laser Doppler Velocimetry and Flowmetry: New Delivery and Detection Techniques", *App. Opt.*, **30**, 2073-2078 (1991)
34. W.F. Van Pelt, W.R. Payne and R.W. Peterson, "A review of selected Bioeffects for Various Spectral Ranges of Light," DHEW Publication (FDA) 20-24, 74-8010
35. see American National Standards Institute, "Safe Use of Lasers": Z-136.1, 1436 Broadway, New York, NY 10018 (1986)

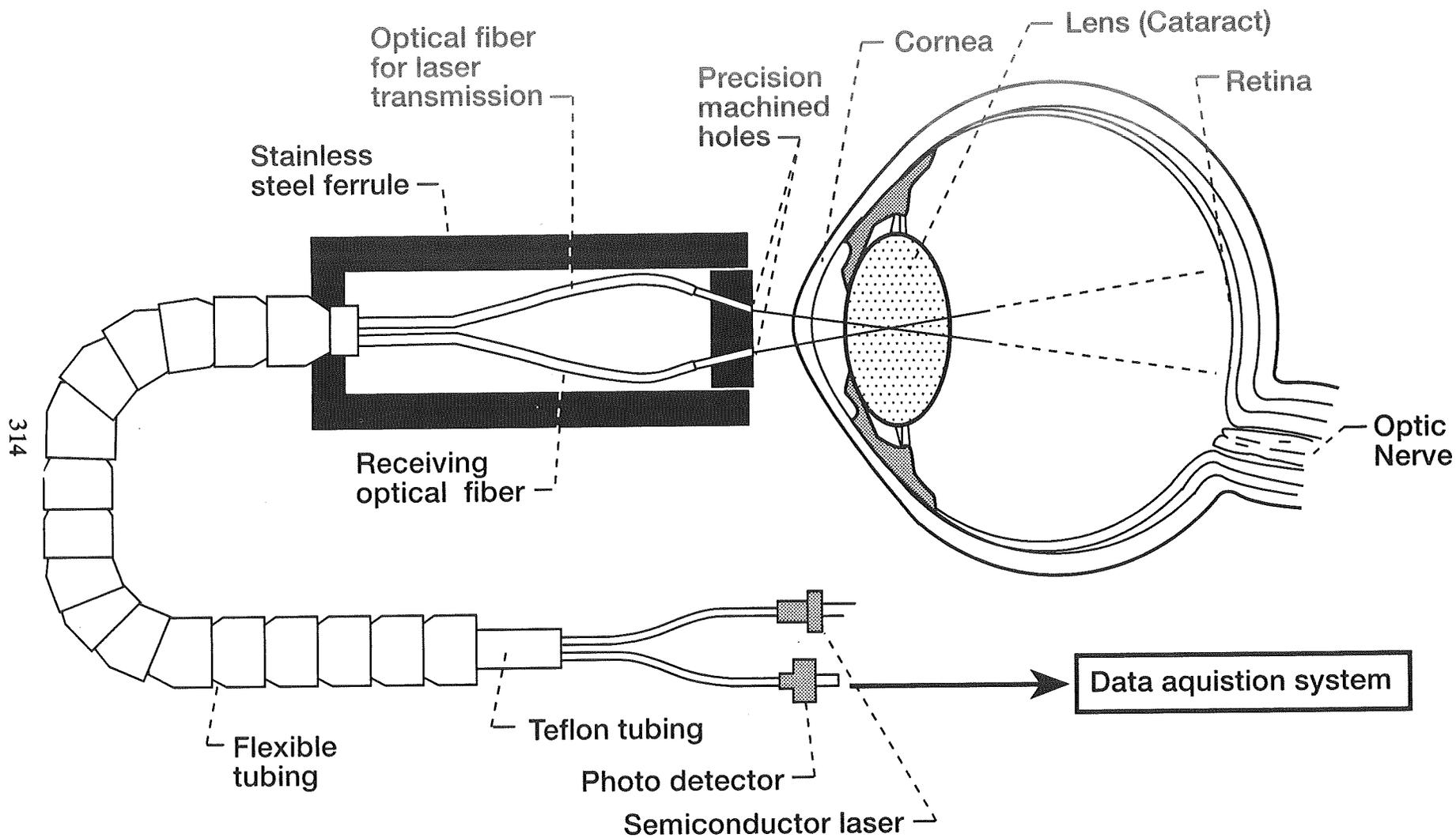


Figure 1: Fiber optic Probe and the Experimental Set up: Schematic of the optical system used for measuring the α crystalline size in an eye lens. The optical design is based upon the dynamic light scattering (DLS) principles.

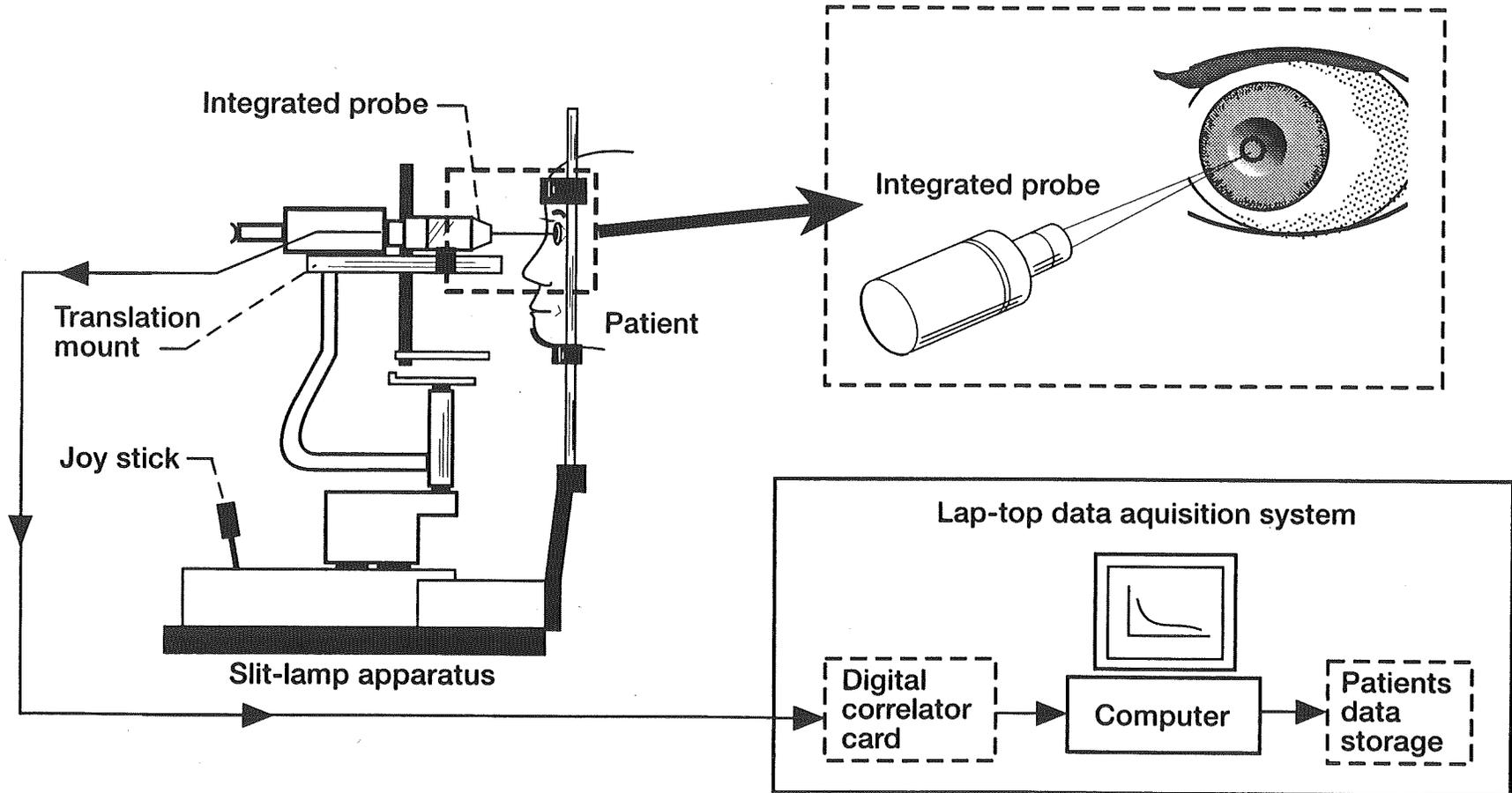


Figure 2: Clinical Set up: The integrated probe is incorporated into a slit-lamp apparatus; commonly used by ophthalmologists for regular eye examinations. The data analysis is performed on a lap-top data computer system, comprising of a digital correlator, a 486 processor, and data storage capability.

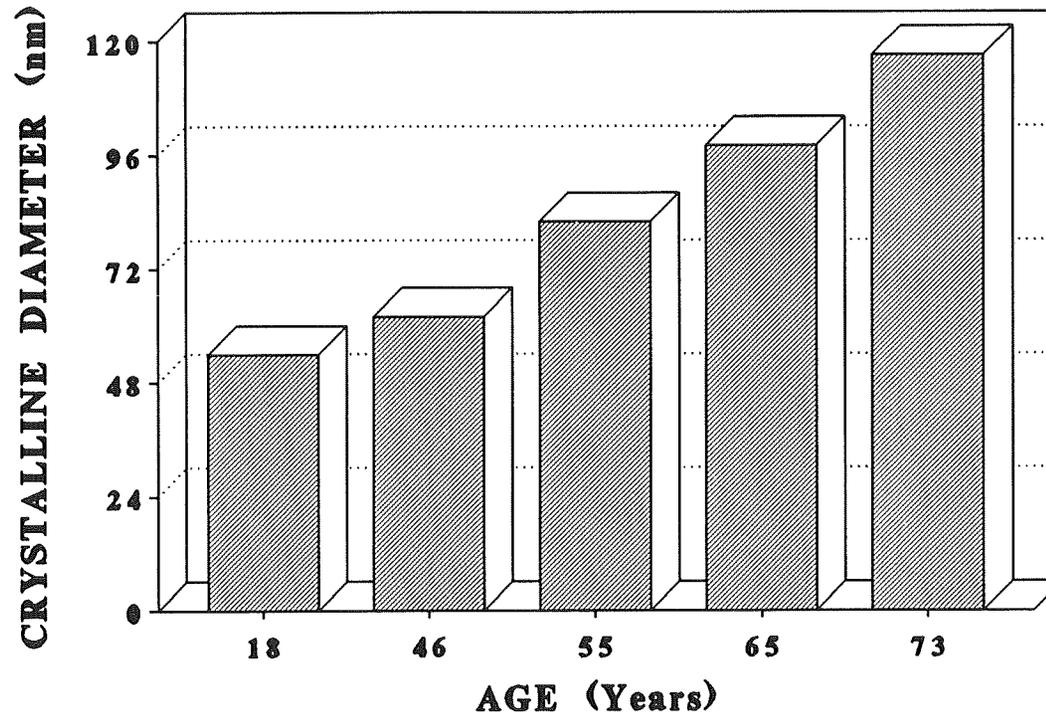


Figure 3: Aging effects (senile cataract) in human eye lenses: The average crystalline sizes are determined using the method of dynamic light scattering. Effective hydrodynamic diameter is computed from Stokes-Einstein relation by using the viscosity of water, and assuming spherical shape for the α crystalline protein molecules in the eye lens.

**MAC TO VAX CONNECTIVITY -
HEARTRATE SPECTRAL ANALYSIS SYSTEM**

Hasan H. Rahman
Monazer Faruque
NASA, Johnson Space Center
GE Government Services
Houston, TX 77058

533-52
150503
p-5

ABSTRACT

The heart rate Spectral Analysis System (SAS) acquires and analyzes in real-time the Space Shuttle on board experiments electrocardiograph (EKG) signals, calculates the heartrate, and applies a Fast Fourier Transformation (FFT) on the heart rate. The system also calculates other statistical parameters such as the 'mean heart rate' over specific time period and a heart rate histogram. This SAS is used by NASA Principal Investigators as a research tool to determine the effects of weightlessness on the human cardiovascular system. This is also used to determine if Lower Body Negative Pressure (LBNP) is an effective countermeasure to the orthostatic intolerance experienced by astronauts upon return to normal gravity. In microgravity, astronauts perform the LBNP experiment within the middeck of the Space Shuttle. This experiment data is downlinked by the orbiter telemetry system, then processed and analyzed in real-time by the integrated Life Sciences Data Acquisition (LSDS) - Spectral Analysis System. The data system is integrated within the framework of two different computer systems, VAX and Macintosh (Mac), using the networking infrastructure to assist the investigators in further understanding the most complex machine on Earth - the human body.

INTRODUCTION

Real-time data analysis has become popular among scientists and researchers as computers become faster and more cost effective. The Spectral Analysis System (SAS) is designed to perform real-time data analysis on microgravity- based life sciences experiments. Over the years, in the life sciences area, scientists have thought of microgravity as a source of physiological difficulties. To study the astronauts adaptation and the after-effect in normal gravity when they return to the earth, the SAS can be a useful tool.

In the past, data were acquired in real-time, but processed and analyzed postflight. The primary focus in real-time was to view the very essential ancillary data during a Space Shuttle mission. The scientist has to retrieve the science data from the NASA data facility to be able to work in their lab and perform further analysis, filtering, digital signal processing and statistical algorithms are needed to study the effects of microgravity and radiation on living organisms.

The SAS is designed to run in real-time on an Apple Macintosh IIfx based desktop computer system. The SAS is capable of acquiring analog EKG samples to perform a realtime FFT output of the heart rate using the off-the-shelf LabVIEW™ 2 software package. The analog signal produced by the LSDS is interfaced into an analog acquisition (NB MIO-16) and Direct Memory Access (DMA) board on the SAS Macintosh. This subsystem acquires data at a rate of 500 samples per second using the LabVIEW™ software and the acquisition hardware. The SAS was used extensively in conjunction with the Life Sciences Data System for analyzing the data in the Science Monitoring Area (SMA) at Johnson Space Center in realtime during the STS 43, STS 44, STS 49, USML and SL-J missions.

The intent of this paper is to introduce the underlying usage of a desktop system like Macintosh incorporating with network interface can help to understand the connectivity issues to retrieve data for real-time analysis.

THE MAC TO VAX CONNECTIVITY

The primary computer systems for the LSDS are VAX computers, which are directly interfaced with the Shuttle downlink telemetry system. The Macintosh computer is used as a backend desktop system for the SMA to aide the investigator in studying the life sciences microgravity-based experiment

data. The VAX and the Mac connectivity is established via the IEEE 802.3 Ethernet, using the DECnet as a protocol for software communication.

The LSDS VAX system acquires, archives and processes the high rate raw downlink multiplexed telemetry data. This data is acquired in realtime and is demultiplexed and extracted by the VAX system for the distribution over the Ethernet network to the desktop system. The Macintosh system uses the 'DECnet for Macintosh' for logical link to access the converted data from the VAX. The extracted data is in a digital form, so further processing is used by the Macintosh computer to convert the data into its original form (which is analog for EKG) to perform the FFT output of the heart rate. Figure 1 shows the overall data flow.

The overall design environment provides a consistent user interface by the Macintosh graphics user interface. The SAS data interchange standards is based on consultative Committee for Space Data Systems (CCSDS) packet standard using the DECnet as a protocol for communication link and access.

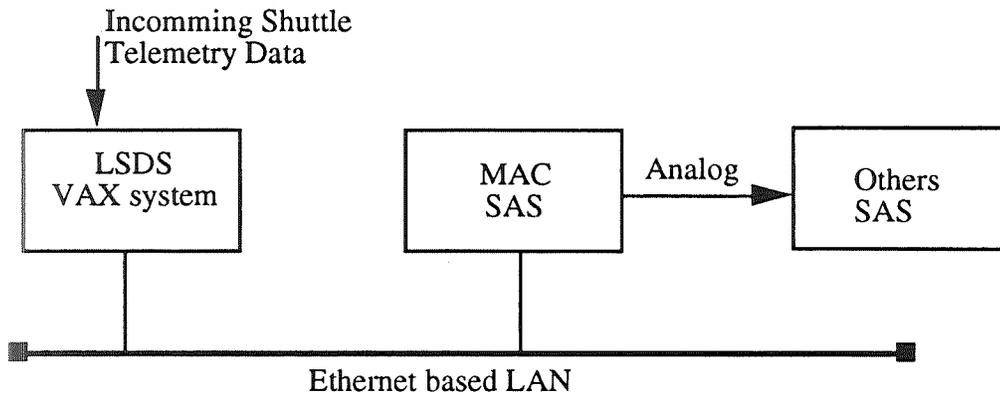


Figure 1: The overall systems data flow

Mac-DECnet Object Declaration

This section provides a sample 'C' program to declare an object for logical link within the Macintosh using DECnet protocol to connect to a VAX system for accessing the network data and command. In this example the Macintosh application waiting for a connection request from the VAX to establish the communication link:

Sample program

```

// set up a DECnet object using DECnet
Boolean      SetupObject()
{
    short error ;
    IOParam param_blk ;
    unsigned char  buffer[3] ;
    Boolean flag = true ;

    //      if (!Control.link_flag)
    //      {
    //          describe the command to setup object
    //          buffer[0] = 2 ;

    //          object number
    //          buffer[1] = Object ;
  
```

```

        buffer[2] = 0 ;
//      DECnet data structure handle
        param_blk.ioPosOffset = (long) Control.Dec_hdl ;
//      DECnet driver reference number, obtained after opening a driver

        param_blk.ioRefNum = Control.ref_num ;
        param_blk.ioCompletion = nil ;

//      number of bytes sent to the driver in this command
        param_blk.ioReqCount = 3 ;

//      command buffer
        param_blk.ioBuffer = &buffer[0] ;

//      declare object by sending message to DECnet
        error = PBWrite((ParmBlkPtr) &param_blk, false) ;

        if (error != noErr)
        {
            ok_dialog(error_box,"DECnet Unable to Setup Object",error, nil) ;
            flag = false ;
        }
    }
    return(flag) ;
}

```

VAX-DECnet Requesting a Connection

This section describes the VAX nontransparent task-to-task DECnet communication setup, using a higher level programming. The 'FORTRAN' example here shows how to create and use a VMS mailbox for receiving network messages, including network status notifications.

```
CALL SYSS$CREMBX(MBX_CHAN,%VAL(MSG_SIZE),%VAL(BUF_SIZE),,,MBX_DESC)
```

```
CALL SYSS$ASSIGN('_NET:',CHAN,,MBX_DESC)
```

The following example provides how to declare a task in VAX enabling to process the Macintosh inbound logical link connection request:

```
!Build a Network Control Block
NCBDESC1(2) = %LOC(NCB)
```

```
CALL SYSS$QIOW(%VAL(CHAN),%VAL(IOS_ACCESS),IOSB,,,NCBDESC1,,,,)
```

The general concept implicit in DECnet-VAX task-to-task communication are covered in Guide to DECnet-VAX Networking Manual by Digital Equipment Corporation.

MACINTOSH SOFTWARE

Once the network data packet is received by the Macintosh system, the data packet is converted into analog samples using the National Instrument's NB-AO-6 board. This sample data is then passed through another National Instrument's NB-MIO-16L board using the DMA features of the National Instrument boards. The above design feature is used because the Macintosh IIfx computer do not provide any DMA - I/O capability. And also, this design allows multiple SAS subsystems for the FFT output from a single network Macintosh connection. In this concept the analog data samples are distributed to other

Macintosh systems, while a single Macintosh is connected to the LSDS VAX system, handling the connectivity and other handshakings.

The analog EKG samples are then passed through a software peak detector to find the peaks for the heart rate calculation. A circular buffer system is designed to implement the acquisition part of the SAS analysis package. Once 256 heart rate samples are collected, an FFT algorithm is run on the array of that heart rate. After the FFT magnitude, the power of different bands and HEART RATE/TIME are calculated and displayed, the next batch of sample data, which is buffered during the intensive calculation of peak detection and FFT algorithm is used to repeat the same process every 10 seconds. Figure 2 shows the data path of the FFT analysis.

The SAS application runs continuously in acquisition mode to acquire EKG signal and calculate the heart rate data from the EKG signal. The output (FFT and the histogram) can be saved as a Pict file onto the disk for further post-test analysis.

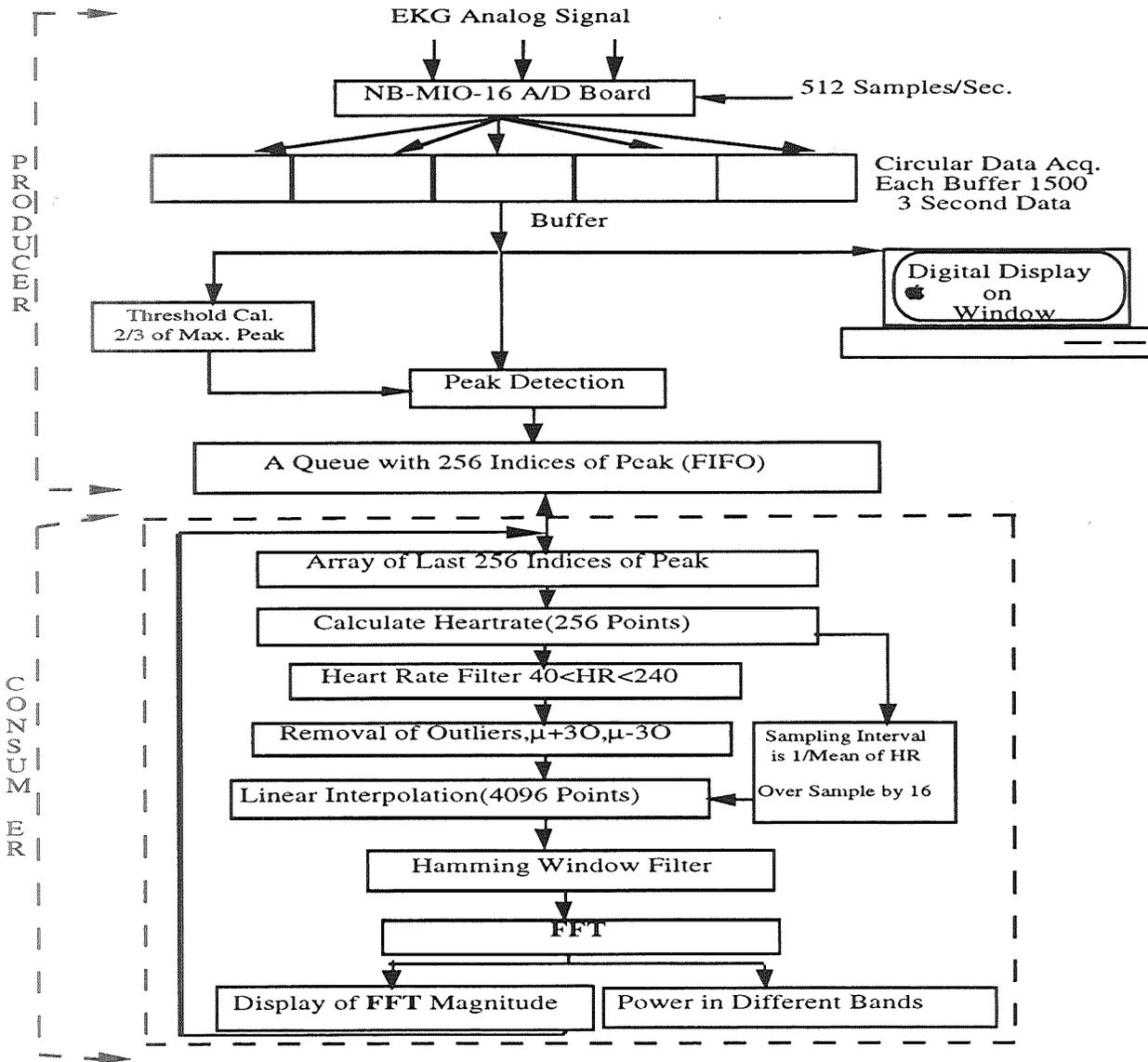


Figure 2. Block Diagram for Data Flow

SOFTWARE REQUIREMENT

LABVIEW 2 Provides the software development environment which is used to program the SAS on a Macintosh platform for data acquisition, manipulation of the acquired data, analyze and display of the acquired experiment data samples.

DECnet for Macintosh provides the network connectivity tools to access the needed data from a VAX system for analysis.

HARDWARE REQUIREMENT

Macintosh IIfx - Provides the fundamental processing platform.

Ethernet Controller Card - Provides the network connectivity.

Digital to Analog Board (NB-AO-6) - Provides the conversion of Digital data into analog samples.

Multi-functional I/O Board (NB-MIO-16L) - Provides the acquisition of analog samples.

Direct Memory Access Board (NB-DMA) - Provides the data transfer from/to I/O board with a minimum interruption and CPU usage.

ACKNOWLEDGEMENT

The Spectral Analysis System described in this paper was developed by Monazer Faruque of the General Electric Government Services, at the NASA Johnson Space Center in Houston, Texas, under a contract with the Life Sciences Project Division.

REFERENCES

1. Digital Equipment Corporation - Guide to DECnet-VAX Networking Manual
2. Digital Equipment Corporation - DEC Pathworks for Macintosh.
3. National Instrument Corporation - LabView 2 Programmers/Users Guide.
4. Robert W. Ramirez, "The FFT Fundamentals and Concepts", 1985 by Tektronix Inc. Englewood Cliffs, New Jersey 07632.

534-53
150504

N93-25595

**THE APPLICATION OF INTEGRATED KNOWLEDGE-BASED
SYSTEMS FOR THE BIOMEDICAL RISK ASSESSMENT
INTELLIGENT NETWORK (BRAIN)**

Karin C. Loftin, Krug Life Sciences, Inc., 1290 Hercules, Suite 120, Houston, Texas, 77058.
Bebe Ly, NASA, Johnson Space Center, Mail Code PT41, Houston, Texas, 77058
Laurie Webster, NASA, Johnson Space Center, Mail Code ER221, Houston, Texas, 77058
James Verlander, Krug Life Sciences, Inc., 1290 Hercules, Suite 120, Houston, Texas, 77058
Gerald R. Taylor, NASA, Johnson Space Center, Mail Code SD5, Houston, Texas, 77058
Gary Riley, NASA, Johnson Space Center, Mail Code PT41, Houston, Texas, 77058
Chris Culbert, NASA, Johnson Space Center, Mail Code PT41, Houston, Texas, 77058
Tina Holden, Lockheed Engineering and Sciences, Corp., Houston, Texas, 77058
Marianne Rudisill, NASA, Johnson Space Center, Mail Code SP34, Houston, Texas, 77058

ABSTRACT

One of NASA's goals for long duration space flight is to maintain acceptable levels of crew health, safety, and performance. One way of meeting this goal is through the Biomedical Risk Assessment Intelligent Network (BRAIN), an integrated network of both human and computer elements. BRAIN will function as an advisor to flight surgeons by assessing the risk of in-flight biomedical problems and recommending appropriate countermeasures. This paper describes the joint effort among various NASA elements to develop BRAIN and an Infectious Disease Risk Assessment (IDRA) prototype. The implementation of this effort addresses the technological aspects of: (1) knowledge acquisition, (2) integration of IDRA components, (3) use of expert systems to automate the biomedical prediction process, (4) development of a user-friendly interface, and (5) integration of the IDRA prototype and Exercise Countermeasures Intelligent System (ExerCISys). Because the C Language, CLIPS (the C Language Integrated Production System), and the X-Window System were portable and easily integrated, they were chosen as the tools for the initial IDRA prototype. The feasibility was tested by developing an IDRA prototype that predicts the individual risk of influenza. The application of knowledge-based systems to risk assessment is of great market value to the medical technology industry.

INTRODUCTION

One of NASA'S primary goals for space flight is to maintain acceptable levels of health, safety, and performance of the crew. To achieve this goal, medical teams have monitored the health of the crew pre-flight, in flight, and post flight throughout the history of manned space programs. During the Skylab missions, in-flight biomedical data were used as a basis for making decisions about the flight duration of successive Skylab missions (21). The medical team had to plan very carefully for a quick turn-around time of sample processing and data analysis between missions. The Skylab Medical Management Group met daily to review the status of the crew. Without computer assistance this activity was very man-hour intensive and undoubtedly increased the cost of the operations.

For extended tours of duty on Space Station Freedom and Lunar/Mars stations, a greater effort will be required to assure nominal crew operations. To achieve this, NASA will monitor crew physiological, psychological, and task performance and administer appropriate countermeasures (17,22,31,39). It may be crucial to assess quickly the individual risk of biomedical problems based on changes in certain physiological, psychological, or environmental indicators to initiate countermeasures (10,19). It is important to predict the impact of the selected countermeasures on crew health, safety, and performance. If more than one change in crew status is observed, it is critical to evaluate each countermeasure relative to the others.

Automation technology is required to support this decisionmaking process. It reduces the volume of data, facilitates data interpretation, and resolves incompatible data. For example, expert or knowledge-based systems can automate the medical diagnostic process (20,28,42). The knowledge that is represented in medical textbooks and/or the expertise of a physician is incorporated into computer software (13). These systems handle a large quantity of related physiological or anatomical data; however, each expert system is developed for only one specific discipline (5,6,7).

Expert systems are commonly implemented as rule-based production systems based on a series of "if...then" reasoning rules (13). The Software Technology Branch at NASA/Johnson Space Center has developed a rule-based production system called CLIPS, the C Language Integrated Production System (9,12,38). CLIPS is being used to automate the prediction process of the IDRA prototype (see below) and BRAIN.

NASA has supported the development of four life sciences expert systems for use on long duration space flight:

- The IDRA prototype assesses the risk of infectious diseases and recommends countermeasures to reduce the risks. The implementation approach and the results of this development are presented in this paper.
- The ExerCISys prescribes an exercise protocol to maintain muscle strength and cardiovascular aerobic capacity in flight.
- The Medical Equipment Computer (MEC) provides decision-support for disease diagnosis and drug therapy in flight.
- The Performance Prediction Model (PPM) assesses the effect of environmental and mission factors on the team performance and predicts its level accordingly. This project is in the early stages of development.

The limitations of these expert systems are that they are independent from each other. They are designed for a single user, and the data are not automatically shared between systems or users. A solution to the problems associated with multiple expert systems is the Biomedical Risk Assessment Intelligent Network (BRAIN). The application of knowledge-based systems or artificial intelligence is a vital component of BRAIN. BRAIN is an integrated network for biomedical risk assessment and management. It provides a composite of multiple experts similar to that generated by the Medical Management Group during the Skylab missions. We hypothesize that BRAIN will reduce the time required for a flight surgeon to arrive at real-time decisions about individual biomedical risk analysis and management. This paper describes an implementation approach to develop a BRAIN prototype and the development of an IDRA prototype to test the feasibility of the approach.

Others outside NASA will benefit from the development of BRAIN. Institutions such as hospitals, medical clinics, boarding schools, military services, nursing homes for the mentally and physically handicapped, and home medical care services are potential users.

IMPLEMENTATION APPROACH

Documentation of Requirements

The BRAIN concept is illustrated as a triangle (Figure 1) with users on the left side and expert systems on the right side. The network still permits each user and system to work independently and interact independently with the flight surgeon. Through BRAIN, each system may access pertinent data from other systems. BRAIN cooperates with the independent expert systems by use of a knowledge base that relates all of them.

The functions of BRAIN are to:

- access IDRA, ExerCISys, MEC, PPM, and other undefined separate systems for pertinent information.
- assess a composite biomedical risk and recommend countermeasures.
- function as a clearing house of information to be shared between systems.
- resolve incompatible information given by other expert systems and derive a composite recommendation for a flight surgeon.

The preliminary software requirements of BRAIN will be documented according to the IEEE Standards Board and the American National Standards Institute (2).

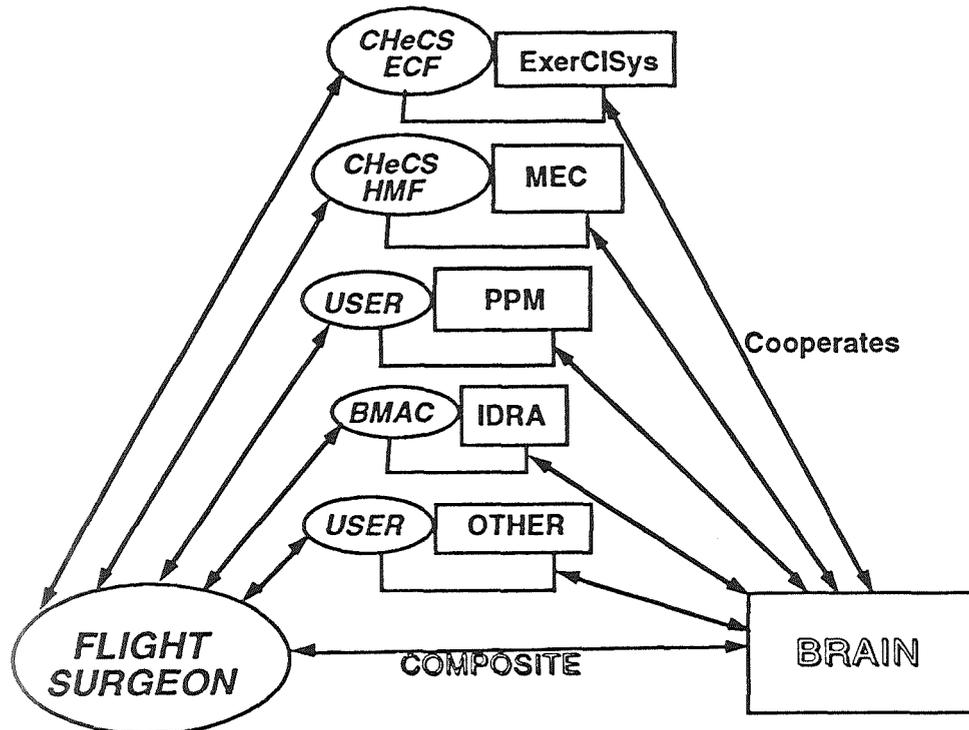


Figure 1. The BRAIN Concept. The users on the left interact verbally with the flight surgeon and mechanically with each independent expert system. BRAIN cooperates with each expert system using a knowledge base that relates all of them. A composite recommendation is then presented to flight surgeon for real-time decisionmaking. (CHeCs-Crew Health Care System, ECF-Exercise Countermeasures Facility, HMF-Health Maintenance Facility, BMAC-Biomedical and Countermeasures)

System Design

Knowledge Definition

A major activity of this project is to develop the knowledge-based system design. This includes the identification of data sources, knowledge definition, knowledge design, and the architecture of the hardware/software environment for BRAIN. The knowledge definition task defines the knowledge requirements of the network and identifies and selects the knowledge sources. The knowledge is acquired, analyzed, and extracted. The knowledge design comprises the knowledge representation, i.e., rules, internal fact structure, detailed control structure, and preliminary user interface (13).

BRAIN receives input from and gives it to PPM, IDRA, MEC, and ExerCISys, or the user, as illustrated in Figure 1. Other data that are required by BRAIN reside in a separate data base, are retrieved as necessary, and are stored in a local or working data base. Unknown data or inaccessible data may be simulated for the version 1.0 development.

The data structure and network configuration of BRAIN must be compatible with IDRA, ExerCISys, MEC, and PPM. These expert systems share related data through BRAIN by accessing the working data base. We will test the feasibility of IDRA, ExerCISys, MEC, and PPM to regularly post data that are required by other systems. A standard protocol will be established for each system to access BRAIN and vice versa.

The knowledge base for BRAIN utilizes and interprets the data, predicts the risk of biomedical problems and recommends the appropriate countermeasures. The information in the data base is extracted from the sources such as:

- *Spaceflight Historical Information*
- *Expert Medical and Science Personnel*
- *Texts, Journal Article and Reviews*
- *Epidemiological Studies of Normal Populations*

The resources available in the medical sciences arena and NASA life sciences groups are explored for the knowledge definition of BRAIN. The relationships among IDRA, ExerCISys, MEC, and PPM are defined by means of workshops, personal consultation and collaboration of existing study results. Experts e.g. flight surgeons or scientists, will be identified and interviewed to model their expertise and to evaluate the demonstration of BRAIN during the developmental stages.

Knowledge Acquisition

Once the knowledge base has been defined for BRAIN, methods will be developed to acquire the specific knowledge. Because a great deal of knowledge has to be acquired for BRAIN, an automated method may be required for that purpose. Several knowledge acquisition tools will be evaluated for consistency and reproducibility in extracting information from human experts and written sources.

The investigative team has access to and experience with several automated knowledge acquisition tools (e.g., Design Alternatives Rationale Tool [DART], Nextra, Task Analysis/Rule Generating Tool [TARGET], Knowledge Acquisition and Representation Tool Kit [KART], and Knowledge Network Organizational Tool [KNOT]).

Nextra operates as a knowledge acquisition front-end tool to an expert system development package called Nexpert Object. Both tools are marketed by Neuron Data from Palo Alto, CA. Nextra allows users to graphically represent entity relationships between the various elements of a subject or domain. DART is another tool that analyzes design alternatives and their associated rationale knowledge. Both Nextra and DART tools can address problems concerned with taxonomies and classification. Both also use repertory-grid knowledge representations. The relationship among the outputs of MEC, ExerCISys, IDRA, and PPM will be examined with these tools in order to derive appropriate rules for biomedical risk assessment. New tools may have to be developed to acquire the appropriate knowledge for BRAIN.

Although various types of expert knowledge exist within the NASA environment, procedural knowledge is prevalent in many areas including the biomedical environment. A procedural analysis tool, TARGET, models a set of actions or procedures associated with a task using a graphical user interface. This tool will be tested to analyze the procedure used by a flight surgeon to solve problems associated with the recommendations given by MEC, ExerCISys, IDRA, and PPM. The specific details will be defined during the knowledge definition phase of the project.

Knowledge Design

A conceptual design of BRAIN is illustrated in Figure 2. Further definition of the knowledge representation and design is delayed until knowledge acquisition is completed. At that time, more will be known about the structure of the knowledge and how it can best be represented.

It is anticipated that the knowledge may be subjected to a software tool called RuleMaster that uses the Iterative Dicotomizer (ID) 3 algorithm. The ID3 algorithm analyses empirical data and derives rules for the knowledge base of BRAIN. Advanced techniques, e.g. CLIPS, will be tested to automate the biomedical prediction process. Other existing and newly-developed tools will be evaluated for their best knowledge representation and design capability.

BRAIN will be designed with a learning capability. It will incorporate, by a feed-back mechanism, the experience of a human expert. The decisions and interpretations of data obtained from actual test cases are acquired automatically in the knowledge base and new rules are induced.

This function is entirely under the control of the appropriate user. But once initiated, it is automatically included in the knowledge base. Tools such as the Automated Reasoning Tool (ART) and Automated Structured Rule Acquisition (ASTRA) are being used to capture the expertise of exercise physiologists for ExerCISys. ART and ASTRA are being evaluated for application to BRAIN.

Knowledge Verification/Validation

Verification and validation of BRAIN is a vital step throughout the life cycle of its development (18). Verification of BRAIN determines that the software is developed according to specifications. The knowledge base will be verified by checking specific details to the level of each rule.

Validation determines that BRAIN performs the functions as specified by the requirements and is usable for field testing (11). Validation of BRAIN will encompass aspects of the validation process such as determining the validation criteria and developing a library of test cases and detailed space flight scenarios that are described in (11) and (27).

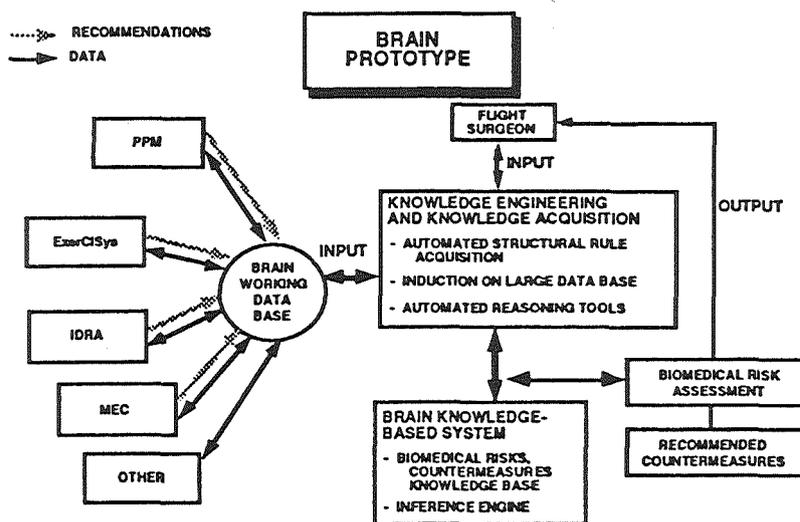


Figure 2. Conceptual Design of BRAIN. The working data base is integrated with IDRA, ExerCISys, MEC, PPM, and additional data bases that contain the facts required by the network, e.g. countermeasures. The expertise of flight surgeons is also captured during the knowledge acquisition process and rules are automatically induced to reflect this expertise. All the rules are stored in the knowledge base and the inference engine executes the appropriate rules for a given working data base.

After the Preliminary Design Review of the project, the detailed design description will be documented. It will specify the logic and content of the knowledge base, the implementation of the system, hardware requirements, the detailed user interface, and the detailed demonstration plan.

The hardware/software environment of BRAIN will be compatible with MEC, PPM, IDRA, ExerCISys, and Space Station Freedom standards to communicate related information. The development environment that is used to create the software may not run on the identical platform as the demonstration version.

User Interface

It is essential for the flight components of BRAIN to have user friendly interfaces. Ease of use may determine whether or not a system is fully utilized. Early prototypes will be developed with prototyping tools to explore the user interface. The user interfaces to BRAIN will be designed in accordance with human factors principles (43) and the Space Station Volume of the Man Systems Integration Standards (NASA STD 3000) document (26). In addition, the user interface code will be portable and compatible with the Space Station Freedom Data Management System. Some of the factors that will be addressed are information grouping, user system dialogs, and information highlighting techniques. Prior to completion of the final BRAIN design, all interfaces will be empirically evaluated using subjects similar to the typical user. Based upon findings of this study, the design of the interfaces will be refined. The final product is BRAIN, version 1.0, that will have been tested and proven to function as an integrated network of the MEC, IDRA, ExerCISys, and PPM prototypes, with a validated, well-designed user interface.

FEASIBILITY TESTING

The IDRA Prototype

The feasibility of using integrated knowledge-based systems for biomedical risk assessment was tested in the IDRA prototype. Because the prevention of infections during manned space flights is important (3,30,35), the IDRA prototype was developed initially to assess the probability of influenza infection.

The epidemiology of and procedures for preventing, diagnosing, and treating influenza are well defined (1,4,40). Epidemiological studies have evaluated the risk factors and their predictive value for influenza in the general population (8,14,25,37,41) and the efficacy of chemotherapeutic prophylaxis (15). Earlier studies investigated the outbreak of influenza in isolated populations, e.g., on an aircraft (29), a ship at sea, (34) and college campuses (24,36). From these sources, we concluded that sufficient information was available to construct a knowledge base about influenza.

The Integration of IDRA and ExerCISys

Studies indicate that exercise has a profound effect on the immune system (23,32,33), sometimes inducing changes similar to those arising from the stress of space flight (16). Therefore, exercise regimen and related physiological data are factors that must be taken into consideration for the risk assessment of infectious diseases and for prescribing an exercise program. This was suggested on the Soviet MIR Space Station when Cosmonaut Gennady Strekalov "caught a cold" following exercise (reported by the Associated Press, October 18, 1990).

The IDRA prototype is compatible with the ExerCISys prototype, and we will integrate IDRA with the ExerCISys as a model for BRAIN. When BRAIN is actually implemented as described in the approach, BRAIN, IDRA, and ExerCISys will be separate systems. However, to test the feasibility of integrating two independent systems initially, the integrated knowledge base will reside in the IDRA prototype. This also provides an opportunity to evaluate the type of information that will be shared and what will remain private between the systems. When the requirements for BRAIN are better defined, the integrated knowledge base will be moved to a separate hardware and software environment, and all systems will be connected to BRAIN through a network communications link.

IDRA Preliminary Results

The knowledge for the IDRA knowledge base was extracted and analyzed from textbooks and journal articles cited above. We identified the critical indicators that predict the probability of influenza. The risk of influenza for an individual is described by general population statistics. It depends on an individual's location, age group, and level of immunity. This information is encoded in a set of 40 rules using CLIPS. Two examples of the rules are in Table I. A subset of these rules incorporates the effect of exercise on the risk of infections. Depending on the individuals condition at any given point in time, a risk of influenza can be assessed based on epidemiological data and the individual's medical record. Once the communication link is completed between the IDRA prototype and the ExerCISys prototype, IDRA will query ExerCISys for the level of fitness of each subject based on aerobic capacity. This information will execute additional rules by IDRA that generate a risk assessment of influenza. The prototype requires further development, validation, and testing.

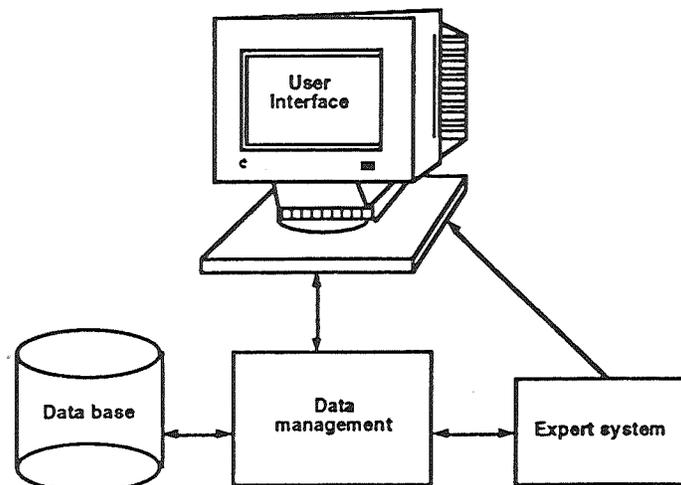


Figure 3. Major Components of the IDRA Prototype

Figure 3 illustrates the major components of the IDRA prototype. A C-based data manager interacts with all the components of the system. It processes information from the data base and from the user interface. The expert

system using CLIPS assesses the probability of influenza. It retrieves the information from the data manager and outputs it to the user interface. For the preliminary user interface, we used the X-Window System. The probability of infection and illness is displayed in the form of text and a graph. All tools are portable and compatible with Space Station Freedom requirements. The preliminary results suggest that an integrated IDRA prototype is feasible and can serve as a model to develop BRAIN.

Table I. Examples of IDRA Rules. The normal state is defined by a., and the effect of moderate exercise is defined by b.

```
a. (defrule normal-state"
  (phase disease-prediction)
  (personal-data (name ?name) (identification ?id)
    (ages ?x&: (or (<?x 18) (> ?x 64)))
    (environment normal) (location ~Houston)
    (nasal-sIgA ?n&: (< ?n 2.75)) (flu-vaccination ~yes))
    (amantadine no) (flu-exposure no) (exercise light))

=>
(update-risk-factor ?id 0.527 0.428); 0.527 is mathematically calculated
(printout t "Subject      :" ?name crlf)
(printout t "Identification : " ?id crlf)
(printout t "age          : " ?x crlf)
(printout t "has a 42.8% chance to get influenza illness due to the age" crlf)
(printout t " and lacking of NASAL sIgA" crlf crlf))
```

```
b. (defrule moderate-exercise-state
  (phase disease-prediction)
  (personal-data (name ?name) (identification ?id)
    (ages ?age) (location ?)
    (environment ~crowded)
    (nasal-sIgA ?n&: (< ?n 2.75)) (flu-vaccination ~yes)
    (amantadine no) (flu-exposure no) (exercise moderate))

=>
(update-risk-factor ?id 0.38 0.28)
(printout t "Subject  :" ?name crlf)
(printout t "Identification : " ?id crlf)
(printout t "age      : " ?age crlf)
(printout t "has only a 28% chance to get influenza illness due to" crlf)
(printout t "moderate exercise" crlf)
(printout t "and lacking of NASAL sIgA" crlf crlf crlf))
```

REFERENCES

1. Atmar, R. L., Greenberg, S. B., Quarles, J. M., Wilson, S. M., Tyler B., Feldman, S., and Couch, R. B.: 1990, Safety and pharmacokinetics of rimantadine small-particle aerosol. *Antimicrobial Agents and Chemotherapy* 34, 2228-2233.
2. Beall, J. H.: 1984, *IEEE Guide to Software Requirements Specifications*. The Institute of Electrical and Electronics Engineers, Inc., New York.
3. Beisel, W. R. and Talbot, J. M. (eds.): 1985, *Research Opportunities on Immunocompetence in Space*. Federation of American Societies For Experimental Biology, Bethesda.

4. Betts, R. F., Douglas, Jr R. G.: 1990, Influenza virus. In: G.L. Mandell, R.G. Douglas, Jr, and J.E. Bennett (eds.), *Infectious Diseases*. Churchill Livingstone, New York, pp. 1306-1325.
5. Blomberg, D.J., Guth, J.L., Fattu, J.M., and Patrick, E.A.: 1986, Evaluation of a new classification system for anemias using Consult Learning System. *Computer Methods and Programs in Biomedicine* 22, 119-125.
6. Bratko, I., Mozetic, I., and Lavrac, N.: 1989, *KARDIO: A Study in Deep and Qualitative Knowledge for Expert Systems*. The MIT Press, Massachusetts.
7. Buchanan, B.G. and Shortliffe, E.H.: 1984, *Rule-Based Expert Systems: The MYCIN Experiments of the Stanford Heuristic Programming Project*. Addison-Wesley Publishing Company, Menlo Park, California.
8. Clover, R.D., Abell, T., Becker, L.A., Crawford, S., and Ramsey, C.N.: 1989, Family functioning and stress as predictors of influenza B infection. *Journal of Family Practice* 28, 535-539.
9. Culbert, C.: 1989, *CLIPS Reference Manual*. National Aeronautics and Space Administration, Houston, JSC-22948.
10. Fry, R. J. M.: 1989, Radiobiological features of the space radiation environment. *Guidance On Radiation Received in Space Activities*. National Council on Radiation Protection and Measurements, Bethesda, MD, pp. 50-144.
11. Geissman, J. R. and Schultz, R. D.: 1991, Verification & validation of expert systems. In: U. Gupta (ed.), *Validation And Verifying Knowledge-Based Systems*. The Institute of Electrical and Electronics Engineers, Inc., Washington, pp. 12-19.
12. Giarrantano, J.: 1989, *CLIPS User's Guide*. National Aeronautics and Space Administration, Johnson Space Center, Houston.
13. Giarrantano, J. and Riley, G.: 1989, *Expert Systems*. PWS-KENT, Boston.
14. Glezen, W.P., Decker, M., Joseph, S.W., and Mercready, Jr., R.G.: 1987, Acute respiratory disease associated with influenza epidemics in Houston, 1981-1983. *Journal of Infectious Diseases* 155, 1119-1126.
15. Glezen, W.P., Grose, N., Haddock, A., and Couch, R.B.: 1989, Chemotherapy and management of respiratory virus infections. *Proceedings of the 16th International Congress of Chemotherapy*. 1989, Israel.
16. Gmunder, F.K., Lorenzi, G., Bechler, B., Joller, P., Muller, J., Ziegler, W.H., and Cogoli, A.: 1988, Effect of long-term physical exercise on lymphocyte reactivity: similarity to spaceflight reactions. *Aviation Space and Environmental Medicine* 59, 146-151.
17. Goldberg, J. et al.: 1987, *A Strategy for Space Biology*. National Academic Press, Washington, D.C.
18. Gupta, U.: 1991, *Validating and Verifying Knowledge-Based Systems*. The Institute of Electrical and Electronics Engineers, Inc., Washington.
19. Haley, R.W., Aber, R.C., and Bennett, J.W.: 1986, Surveillance of nosocomial infections. *Hospital Infections, 2nd Edition*. Little, Brown & Co., Boston, pp. 51-71.
20. Hand, D.J.: 1987, Artificial intelligence and medicine. *Journal of the Royal Society of Medicine* 80, 563-565.
21. Johnston, R.S.: 1977, Skylab medical program overview. In: R.S. Johnston, and L.F. Dietlein (eds.), *Biomedical Results From Skylab*. National Aeronautics and Space Administration, Washington, D.C., pp. 1-19.
22. Kanas, N.: 1985, Psychosocial factors affecting simulated and actual space missions. *Aviation Space and Environmental Medicine* 56, 806-811.

23. Keast, D., Cameron, K., and Morton, A.R.: 1988, Exercise and the Immune Response. *Sports Medicine* 5, 248-267.
24. Layde, P.M., Engelberg, A.L., Dobbs, H.I., Curtis, A.C., Craven, R.B., Graitcer, P.L., Sedmak, G.V., Erickson, J.D., and Noble, G.R.: 1980, Outbreak of influenza A/USSR/77 at Marquette University. *Journal Infectious Diseases* 142, 347-352.
25. Longini, Jr., M., Koopman, J.S., Haber, M., and Cotsonis, G.A.: 1988, Statistical inference for infectious diseases. Risk-specific household and community transmission parameters. *American Journal of Epidemiology* 128, 845-859.
26. Man-Systems Integration Standard, Vol. IV, NASA STD-3000, Rev. A, June 1991.
27. Marcot, B.: 1991, Testing your knowledge base. In: U. Gupta (ed.), *Validating And Verifying Knowledge-Based Systems*. The Institute of Electrical and Electronics Engineers, Inc., Washington, pp. 188-199.
28. Miller, P.L.: 1988, Artificial intelligence in medicine: An emerging discipline. In: P.L. Miller (ed.), *Selected Topics in Medical Artificial Intelligence*. Springer-Verlag, New York, pp. 11-24.
29. Moser, M.R., Bender, T.R., Margolis, H.S., Noble, G.R., Kendal, A.P., and Ritter, D.G.: 1979, An outbreak of influenza aboard a commercial airliner. *Journal of Epidemiology* 110, 1-6.
30. Nicogossian, A.E. and Garshnek, V.: 1989, Historical perspectives. In: A.E. Nicogossian (ed.), *Space Physiology and Medicine*. Lea and Febiger, Philadelphia, pp. 17-29.
31. Nicogossian, A.E. (ed.): 1989, *Space Physiology and Medicine*. Lea and Febiger, Philadelphia.
32. Nieman, D.C., Johanssen, L.M., and Lee, J.W.: 1989, Infectious episodes in runners before and after a roadrace. *Journal of Sports Medicine and Physical Fitness* 29, 289-296.
33. Nieman, D.C. and Nehlsen-Cannarella, S.L.: 1991, The immune system. In: C.B. Rians (ed.), *Principles and Practice of Sports Medicine*. Human Kinetics Publishers, Champaign, IL, pp. 1-49.
34. Olson, J.G., Ksiazek, T.G., Irving, G.S., and Rendin, R.W.: 1979, An explosive outbreak of influenza caused by a/USSR/77-like virus on a United States naval ship. *Military Medicine* 144, 743-745.
35. Pierson, D.L.: 1986, *Space Station Infectious Disease Risks*. National Aeronautics and Space Administration, Johnson Space Center, Houston, JSC-32104.
36. Pons, V.G., Canter, J., and Dolin, R.: 1980, Influenza A/USSR/77 (H1N1) on a university campus. *American Journal of Epidemiology* 111, 23-30.
37. Reuman, P.D., Bernstein, D.I., Keely, S.P., Sherwood, J.R., Young, E.C., and Schiff, G.M.: 1990, Influenza-specific ELISA IgA and IgG predict severity of influenza disease in subjects prescreened with hemagglutination inhibition. *Antiviral Res* 13, 103-110.
38. Riley, G.: 1989, *CLIPS Architecture Manual*. National Aeronautics and Space Administration, Houston, JSC-23047.
39. Robbins, F. et al.: 1988, *Exploring the Living Universe/A Strategy for Space Life Sciences*. National Aeronautics and Space Administration, Washington, D.C.
40. Taylor, R., Nemaia, H., Tukuitonga, C., Kennett, M., White, J., Rodger, S., Levy S, and Gust I: 1985, An epidemic of influenza in the population of Niue. *Journal of Medical Virology* 16, 127-136.
41. U.S. Department of Health and Human Services: 1989, *Vital and Health Statistics: Current Estimates From the National Health Interview Survey, 1988*. DHHS Publication, Hyattsville, Maryland, (PHS) 89-1501.

42. Winkel, P.: 1989, The application of expert systems in the clinical laboratory. *Clinical Chemistry* 35, 1595-1600.
43. Woodson, W.E.: 1981, *Human Factors Design Handbook*. McGraw-Hill Book Co., New York.



omit

**ENERGY AND ENVIRONMENT PART 2:
ENERGY INNOVATIONS**

PRECEDING PAGE BLANK NOT FILMED

333.

332
~~INTENTIONALLY BLANK~~

**SOLID-STATE ISOTOPIC POWER SOURCE
FOR COMPUTER MEMORY CHIPS**

535-33
150505

PAUL M. BROWN, Ph.D.
IsoGen Radioisotopic Research Laboratory
315 S. McLoughlin Blvd.
Oregon City, OR 97045
(503) 656-4419

N 93-25596

ABSTRACT

Recent developments in materials technology now make it possible to fabricate nonthermal thin-film radioisotopic energy converters (REC) with a specific power of 24 W/kg and a 10 year working life at 5 to 10 watts. This creates applications never before possible, such as placing the power supply directly on integrated circuit chips. The efficiency of the REC is about 25% which is two to three times greater than the 6 to 8% capabilities of current thermoelectric systems. Radioisotopic energy converters have the potential to meet many future space power requirements for a wide variety of applications with less mass, better efficiency, and less total area than other power conversion options. These benefits result in significant dollar savings over the projected mission lifetime.

INTRODUCTION

Placing miniature power supplies right where they're needed on integrated circuit chips is a quick and efficient way of getting electrical power to a circuit's microscopic components. This goal now seems within reach with the construction of the tiny, thin-film radioisotopic energy converter or REC that generates electricity.

Traditional lead-acid batteries are made from plates of lead and lead peroxide immersed in a sulfuric acid solution, called an electrolyte. Charged atoms, called ions, flow through the acid within the battery from one plate to another, setting up an electrical imbalance that causes electrons to flow through wires attached to the battery terminals and completing the circuit. These batteries are heavy and corrosive.

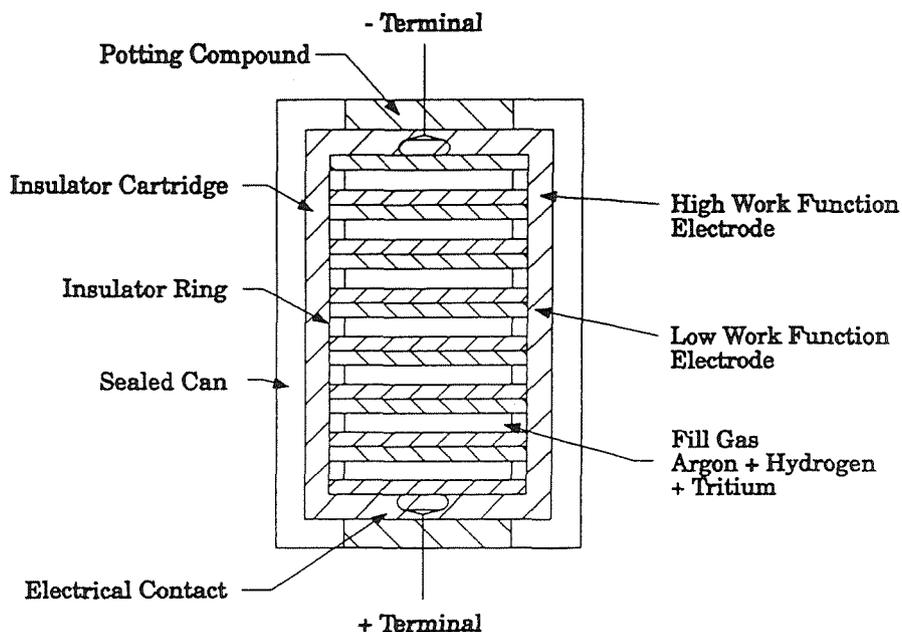


FIGURE 1. Contact Potential Gas Ionization Type.

334
MICROFILMED BY

Interest in non-chemical space power systems has been rekindled with the recent ambitious manned space missions planned by NASA's Space Exploration Initiative. Deep space missions will require more energy than is practically available from chemical systems.

The operation of a contact potential cell was first demonstrated by Kramer in 1924 (Brown 1992a) while new developments have been the subject of much research (Brown 1992b). When a gas is irradiated by beta particles, ionization results. If this occurs in an electric field, such as produced by the contact potential difference of two dissimilar metallic electrodes (Figure 1), which are electrically connected through an external conductor, the ions will migrate to the electrode of opposite charge and conduction through the gas takes place. The electrical energy delivered to any external load is a result of the ionizing energy of the beta particles. The contact potential difference enables the neutralization of ions and the transfer of their charge to the circuit to take place. Efficiency is limited to 1% by charge recombination, leakage current, and space-charge effects but long life of 20 years is easily achieved.

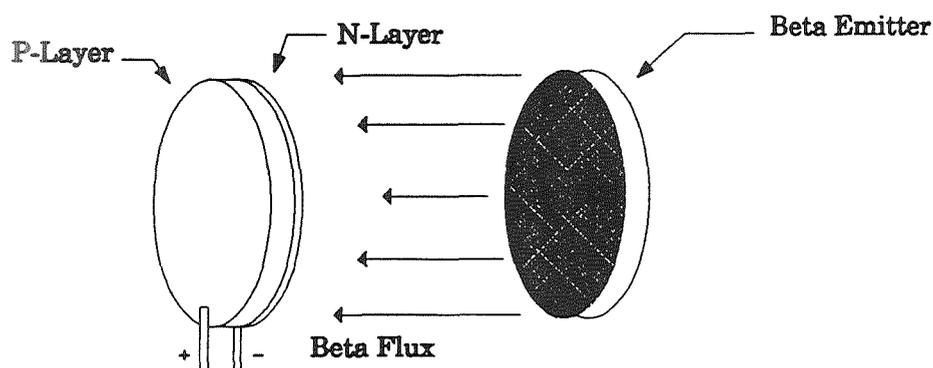


FIGURE 2. Basic Approach to Betavoltaic Energy Conversion Consisting of a Source Coupled to a Semiconductive Junction Device.

The first work on p-n junction converters was reported by Rappaport in 1954 (Linder 1956), since then betavoltaic p-n junction batteries have been used only as a source of power in the heart pacemaker, where power requirements are on the order of micro-watts (Olsen 1972). Figure 2 shows a straightforward approach to a betavoltaic p-n junction energy converter; a planar source of beta emitting material is coupled to a planar p-n junction device. The n-face and p-face are respectively the negative and positive terminals of the single cell power source. Utilizing the internal electric field of the junction, the radiation energy is converted directly into electrical energy, in much the same manner as a solar cell converts photons of light into electricity (Olsen 1974). Efficiency is good on the order of 25% and fairly reliable using low power isotopes in the microwatt range but radiation damage to the junction seriously limits the working life of this type of device, when built with any power above one milliwatt, to only a few days.

I should also mention the most widely used type of nuclear battery which is called the radioisotopic thermoelectric generator (RTG). This type of device utilizes the heat produced by large amounts of radioactive material with a thermo-couple to generate electricity. This type of device is limited to only 6 to 8% efficiency and requires high operating temperatures in addition to the fact that much shielding is necessary because of the type and amounts of radioactive material used. However, these devices are reliable and have been used repeatedly by NASA as well as other agencies.

THE RADIOISOTOPIC ENERGY CONVERTER (REC)

The nonthermal thin-film isotopic energy converter is a hybrid contact-potential/betavoltaic device. It is made up of alternating thin-film layers of low work-function metal (Figure 3), isotope laced semiconductor media, and a high work function metal. Beta particles emitted from the source traverse the media losing energy and creating electron-hole pairs. Those carriers within a diffusion length of the junction will be swept across the gap contributing a current. A United States Patent was granted on this innovative approach to radioisotopic decay energy conversion February 11, 1992, USP #5,087,533.

The embodiment shown here in figure 3 consists of a thin film layer of cesium metal on which a second layer of selenium is deposited. Other materials will also work in place of the selenium. This selenium layer is loaded with tritium for this configuration although any suitable isotope will work. The third or final layer is of platinum metal which yields a composite foil 0.001 mm thick and will deliver 4.4 volts continuously.

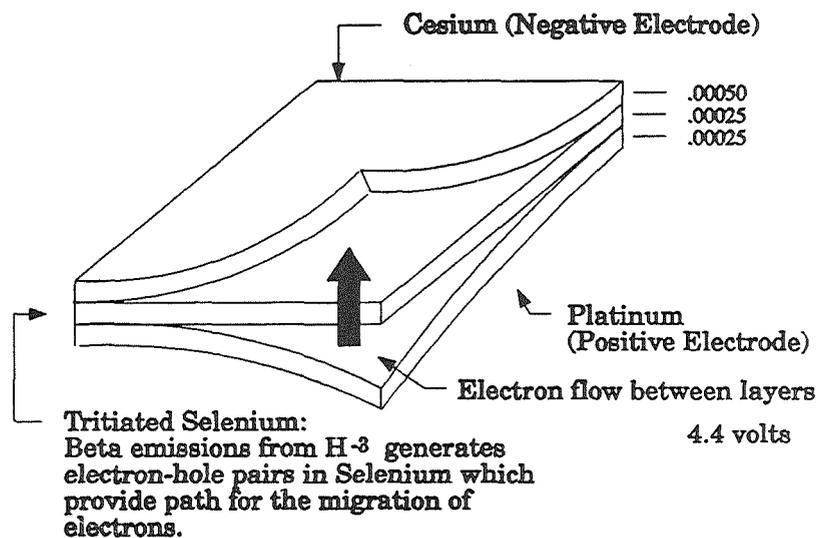


FIGURE 3. Radioisotope Electric Converter (Paper-Thin Layers are Shown Peeled Back for Illustration).

Just as a solar cell converts light energy directly into electrical energy, the REC converts radioactive decay energy directly into electricity. The new technology of the REC works on the same principle as traditional batteries only the ions are not generated by chemical reactions. Rather, the acid solution or electrolyte is replaced by a semiconductor medium that is ionized by the absorption of radioactive decay particles emitted from a radioisotope homogeneously dispersed throughout the semiconductor itself. That is to say, the semiconductor is doped with the radioisotope. The semiconductor medium is a thin-film and conducts ions between metal electrodes typically made of cesium and platinum. Versions of the REC develop a potential difference of 1 to 4.5 volts DC between the electrodes (depending on the type of electrodes used) from a single cell and a few milliwatts of power per square centimeter.

One embodiment of the REC yields 4.4 volts DC from a foil only 0.001 mm thick while another embodiment produces 9 volt DC from a foil .045 mm thick. The current is determined by the isotope used and the surface area of the foil. This foil may be rolled into tubes, spirals or be cut into virtually any shape. Calculations show a specific power of 24 watts per kilogram with a 10 year working life is achievable for a 10 watt device (NASA 1991). Based on recent experimentation the efficiency is estimated to be on the order of 25% thermal to electric, which is three times greater than the 8% capability of current thermoelectric systems.

The REC is superior to betavoltaic devices of the p-n junction type because the REC design does not subject the voltage mechanism to radiation damage; superior to the contact potential difference converter because without a gas the REC is not susceptible to space charge effects while ion recombination is limited by the narrow

media gap; and superior to the RTG by operating at low temperatures with greater efficiency, lower weight, and reduced health dangers.

The possibility of using relatively simple manufacturing techniques to fabricate these miniature, thin-film RECs at low cost suggests a variety of applications. Here at IsoGen, we are investigating how to deposit these devices directly onto integrated circuits to provide onsite power for computer memory chips.

In concept, the REC offers several advantages over batteries made with liquid electrolytes. The REC will last much longer due to the fact that radioisotopic energy is several orders of magnitude greater than chemical energy. They also operate over a broader temperature range than batteries with electrolytes made from liquids, which may freeze in extreme cold or boil in extreme heat. Some batteries on the market use a gel or a solution suspended in plastic as their electrolyte, but none is totally solid state. The REC may be the first totally solid state battery.

"Smart" credit cards equipped with a tiny computer chip to keep track of a user's banking records as well as such other information as medical history are being tested in Europe and Japan. The REC could run smart credit cards, cardiac pacemakers, possibly even artificial hearts or other electronics that do not require large amounts of power. New applications will be developed around this new long life portable power source that we can only dream of today.

Beta-Conductivity

Beta-conductivity is a phenomenon evidenced by the increase in electrical conductivity of a material after the absorption of ionizing radiation. This effect is attributed to the increased number of free electrons generated by absorption of the ionizing radiation. Of course, recombination of the charge carriers does occur, however, during the time these carriers are free in the material, the conductivity can be greatly enhanced. We minimize recombination by using very thin layers of the ionized material in the REC.

For a material in which one type of carrier predominates, for example electrons, the change in conductivity with irradiation can be expressed as (Brown 1990b):

$$\Delta\sigma = \Delta n e \mu + e n \Delta\mu \quad (1)$$

where $\Delta\sigma$ is conductivity change,
 Δn is change in free carrier density,
 e is electronic charge,
 μ is carrier mobility, and
 $\Delta\mu$ is change in carrier mobility.

"Beta-conductivity gain" may be defined as the number of inter-electrode transits that can be made by an electron in a conductor before the beta-generated "hole" is eliminated by recombination. For this case where one type of carrier predominates, the gain is expressed as (Brown 1990b):

$$G = \chi \mu^{V\lambda-2} \quad (2)$$

where χ is carrier lifetime,
 V is applied voltage, and
 λ is spacing between electrodes.

Principles of Operation

It is presumed that free electrons are retained in a metal by the forces of attraction between two unlike charges, since there must be as many positively charged atoms as there are free electrons. An amount of work must be performed to remove an electron from the surface of any metal called the work function of the material, and of course, this value may be calculated. But for this discussion it is sufficient to say that the work function of most elements, including the metals have been measured and tabulated in suitable reference books. Typically, work functions range from 1 eV (easiest to remove) to 5.5 eV (most difficult) for metals.

When two metals of different work functions, say cesium and platinum, are brought into electrical contact (by direct contact or through a conductor) an electrical potential difference is established as discovered by Volta in 1797. This contact potential difference is the result of the difference between the work functions of the two metals. The sign is determined by the fact that the metal with the largest work function gives up electrons the least readily and becomes the more electronegative.

An amount of work must be performed on a neutral atom to remove electrons (ionize the atom). This work manifests itself as increased potential energy and may be utilized to do work before allowing the electron and ion to recombine.

The Betavoltaic Effect may simply be defined as the conversion of ionizing radiation to electrical energy by a material or combination of materials. Radiation that is absorbed in the vicinity of any potential barrier will generate separate electron-hole pairs which in turn flow in an electric circuit due to the Voltaic Effect. Of course this occurs to a varying degree in different materials and geometries.

Figure 4 is a representation of a basic betavoltaic converter and demonstrates the characteristics of the REC. Electrode A has a positive potential while electrode B is negative with the potential difference provided by the Contact Potential Difference of the dissimilar metals. An electric field exists between the electrodes, a zone we will call the junction. The junction between the two electrodes is thus comprised of a suitably ionizable medium exposed to decay particles emitted from a radioactive source.

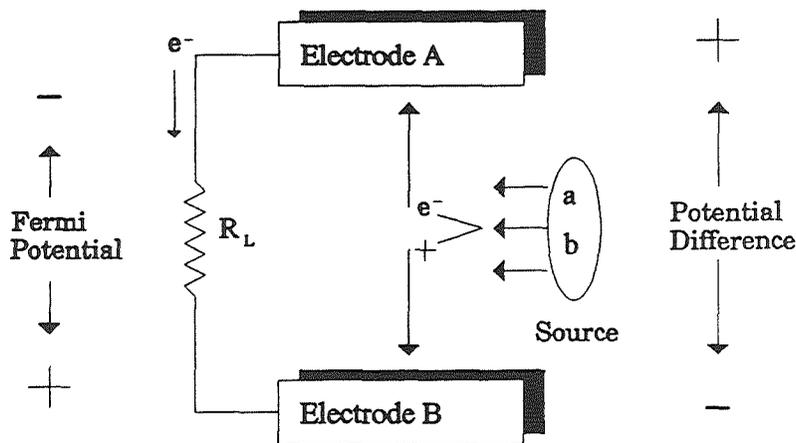


FIGURE 4. Betavoltaic Converter.

In general, the introduction of ions from any source into an electric field will generate electricity in accordance with well known physical principles and may be satisfactorily explained in terms commonly associated with the Volta Effect. The energy contributed to such a circuit does not come from the ions themselves but rather from the work done on the circuit to generate the ions, known as the ionization potential of that particular material.

Neither the electric field, the electrodes or the medium between the electrodes contribute any energy in the Voltaic Effect. The energy is contributed by the ion generator whether this mechanism is chemical, electromagnetic or nuclear is irrelevant. Modeling for the REC has been generated and published (Brown 1990a).

CONCLUSION

The REC is a current generating device including a pair of electrodes of electro-chemically dissimilar materials separated by a space filled with a solid medium having a relatively high dielectric constant and a relatively low ionization potential. Current flows in an external circuit coupling the two electrodes together when the solid medium is ionized. A suitable radioactive material is mixed or dispersed in the solid medium to provide the ionizing flux.

The REC could be a less toxic, longer lasting alternative to conventional chemical batteries, and most likely to compete with button cell batteries to power small electronic items. Such a device could have a broad range of applications, from portable power supplies for integrated circuits to information processing, and nonpolluting generators in portable phones. The REC's lightweight, durability, and lack of corrosiveness also make it a natural candidate for space flight.

Market surveys have been conducted by the nuclear industry in the past and the conclusion has been that there is a need for a long-life radioisotope nuclear battery. Economic studies indicate that the REC could be economically competitive with chemical batteries for applications requiring lifetimes of over two years at remote locations where the expense of charging or changing batteries is significant. Applications where the inaccessibility after implantation is a consideration that leads to the selection of an REC due to its superior reliability and life.

ACKNOWLEDGMENTS

This work was performed at the IsoGen Incorporated Radioisotopic Research Laboratory under internal funding support.

REFERENCES

- Brown, P. (1990a) "Radioisotopic Generators: Beta Voltaic Effect Modeling," *Raum & Zeit Magazine*, 2(1):81-82.
- Brown, P. (1990b) "The Beta Voltaic Effect Applied to Radioisotopic Power Generation," American Nuclear Society Annual Meeting, Nashville, TN.
- Brown, P. (1992a) "Isotopic Electric Generators: Conception to Application," IsoGen, Inc., Portland, OR.
- Brown, P. (1992b) "The Contact Potential Isotopic Generator," *Raum & Zeit Magazine* 3 (3):72-73.
- Linder, E. (1956) "Direct Conversion of Radiation into Electrical Energy," *Proceedings of the International Conference on the Peaceful Uses of Atomic Energy*, 15:283-290, United Nations, NY.
- NASA (1991) "Betavoltaics of Increased Power," *NASA Tech Briefs*, 5(8):32
- Olsen, L. (1972) "Betavoltaic Energy Conversion," *Energy Conversion* 13:117
- Olsen, L. (1974) "Advanced Betavoltaic Power Sources," 9th Intersociety Energy Conversion Engineering Conference

**PHOTOVOLTAIC POWER WITHOUT BATTERIES
FOR CONTINUOUS CATHODIC PROTECTION**

**W. W. Muehl, Sr.
Department of the Navy, Coastal Systems Station (COASTSYSTA)
Dahlgren Division, Panama City, FL 32407-7001**

ABSTRACT

The objective of this project was to successfully demonstrate that renewable energy can efficiently and economically replace dedicated non-renewable power sources.

The COASTSYSTA designed, installed, and started up on 20 January 1990, a state-of-the-art photovoltaic powered impressed current cathodic protection system (PVCPSYS) not requiring any auxiliary/battery backup power for steel and iron submerged structures. The PVCPSYS installed on 775' of steel sheet piling of a Navy bulkhead is continuing to provide complete, continuous corrosion protection well documented by COASTSYSTA and verified on-site by the U.S. Army Corps of Engineers.

The PVCPSYS uses only renewable energy and is environmentally clean. A patent is pending on the new technology. Other possible PVCPSYS applications are mothballed ships, docks, dams, locks, bridges, marinas, and pipelines.

The Department of Defense Photovoltaic Review Committee and Sandia National Laboratories consider this successful and cost effective system a major advance in the application of photovoltaics.

The objective of this project was to successfully demonstrate that renewable energy can efficiently and economically replace or be used instead of continuous non-renewable power sources. An opportunity to clearly show that photovoltaic power is practical was the result of a recommendation to provide cathodic protection to the Naval Diving and Salvage Training Center bulkhead.

The COASTSYSTA in Panama City, Florida, has broken new ground in the application of solar energy for cathodic protection. Photovoltaic arrays without battery backup have been connected to the 775 foot-long steel sheet piling of a dock bulkhead via a cathodic protection system, to prevent corrosion on that steel structure in a salt water environment.

Cathodic protection, as the name signifies, is the process by which, in the COASTSYSTA impressed current type application, the entire steel sheet piling is transformed into a cathode via a series of anodes mounted in PVC standoff racks, in the water, next to the piling. When direct current (DC) energy is applied to the anodes and sufficient electrical potential is attained by current flow from the anodes via an electrolyte (seawater) to the piling, the corrosion is transferred to the anodes, preventing piling corrosion.

Mr. Wally Muehl, Publics Works Engineer at the Coastal Systems Station, was evaluating power sources to protect the Naval Diving and Salvage Training Center bulkhead when he focused on photovoltaics. Although there are 10 other impressed current cathodic protection systems installed on the docks, all are powered by a continuous power source with the current rectified to DC.

The Naval Diving and Salvage Training Center is in a separate location from these docks, and it was determined that power was not readily available and would be expensive to provide rectifiers on the dock due to the dock configuration. Rectifiers would also pose a safety hazard on the dock that is regularly used for diver and salvage training. This bulkhead was 12-years old and other than the initial coating, received no corrosion protection.

Mr. Muehl developed a state-of-the-art solar powered system for impressed current cathodic protection of submerged steel and iron type structures without requiring any battery backup power. Battery backup power is considered costly and an environmental problem. To date, all impressed current systems require a continuous DC power supply in order to provide cathodic protection.

The COASTSYSTA photovoltaic power system is a fixed-axis system which is suitable for the Panama City latitude of 30°10'N, 85°22'W. The tilt of the arrays were set at latitude instead of +15 degrees in January 1990, and have not been changed. This is a good indication that other areas with good distribution, but lower insolation levels, would be excellent prospects for a similar type of photovoltaic powered system. For higher latitudes, there are several other options to improve system performance without battery backup. These include one-axis East-West tracking, two-axis North-South, East-West tracking, or simply adding a module or two to meet the additional current requirements.

As engineer in charge, Mr. Muehl, who designed, prepared the specifications, and monitored the installation, also had two other problems that had to be considered and resolved in order to install a impressed current cathodic protection system. The first problem was ensuring that the steel piling had electrical continuity. Another problem was providing sufficient impression of current "carry over" to overcome a 155-foot section of piling that had to be bypassed, and provide cathodic protection, without anode placement in the area having a water depth of 27 feet, where diving takes place. Both problems were overcome in the design.

To facilitate the use of a photovoltaic powered cathodic protection systems without battery backup, the steel sheet pilings were provided an initial one-time only preconditioning polarization for a predetermined continuous time period to the extent that these pilings were initially polarized to a high negative potential by a temporary DC power source. It is to be noted that the evolution of a protective hydrogen film is merely a beneficial by-product of the preconditioning polarization at the higher negative potentials. The initial DC power for polarization can be provided by a DC power source such as a portable motor driven DC generator or a portable motor driven DC welder.

The COASTSYSTA photovoltaic powered cathodic protection system tests performed and other data obtained, provide a further explanation that the anode-seawater-cathode piling structure acts like a battery and when allowed to rest, the polarity level recovers and is electrochemical in nature. An electrochemical lead-acid battery, for example, can recover charge if allowed to rest after serving a load. The electrochemical reaction reverses slightly when the load is disconnected, however, a capacitor without an external current source cannot recover by simply removing the load. It is believed that the one-time only initial preconditioning polarization of the structure embeds atomic hydrogen which can also migrate and diffuse in the structure. This system delays the decay of the negative potential and permits the photovoltaic

arrays to provide a "trickle" charge, allowing the system to easily provide complete continuous cathodic corrosion protection including cloudy, overcast, rainy and nighttime conditions without the necessity for DC power backup such as batteries.

In summary, the foregoing novel method and system of a one-time-only preconditioning or prepolarizing the structure prior to energizing the PV solar array on-line with the system, provides a relatively higher negative potential that has a slow rate of decay. This permits the use of regulated PV solar energy with excess available power, and without any backup power, to easily provide complete continuous corrosion protection, including cloudy, overcast, rainy and nighttime conditions, with excellent polarization levels and improving with time. An analogy may be that the steel structure becomes very effectively polarized, and will remain so by the variable DC charge effect provided by the simple solar array system, much like a piece of steel or iron can become magnetized by the application of a DC electrical current.

The installation, start up, and continuing operation, including underwater inspections, are well documented to date by the Coastal Systems Station and verified on site, during the day and at nighttime by the U. S. Army Corps of Engineers, Construction Engineering Research Laboratory, Naval Energy Program Office and members of the Department of Defense (DoD) Photovoltaic Review Committee. The average amount of available sunshine for the three weeks prior to these organizations visit, per data provided by the National Weather Service, averaged 24%.

This system has been in operation almost 2 1/2 years without requiring any maintenance or adjustment. A patent is pending on the new technology. Other possible applications are mothballed ships, docks, dams, locks, bridges, marinas, and pipelines.

The Department of Defense Photovoltaic Review Committee and Sandia National Laboratories consider this successful and cost effective system a major advance for the application of photovoltaics.

**PHOTOVOLTAIC POWER WITHOUT BATTERIES
FOR CONTINUOUS CATHODIC PROTECTION**

W. W. Muehl, Sr.

**Department of the Navy, Coastal Systems Station (COASTSYSTA)
Dahlgren Division, Panama City, FL 32407-7001**

With reference to this manuscript, it is with pleasure that I acknowledge the "helpful cooperation and information" received from the following personnel:

- Dr. Michael G. Thomas and Mr. Terry Schuyler - Senior Members, Technical Staff, Photovoltaic Research Dept., Sandia National Laboratories, Albuquerque, NM
 - Mr. James F. Jenkins, P.E., Corrosion & Metallurgical Engineer, Naval Civil Engineering Laboratory, Port Hueneme, CA
 - * Mr. L. E. Humble, Photovoltaic Programs, Energy Program Office, Naval Weapons Center, China Lake, CA
 - * Mr. Roch A. Ducey, Principal Investigator, U.S. Army Construction Engineering Research Laboratory, Champaign, IL
 - Mr. Thomas F. Lewicki, P.E., Facilities Corrosion Program Manager, HQ Air Force Civil Engineering Support Agency, Tyndall Air Force Base, FL
 - Navy Divers & Dive Locker, Coastal Systems Station, Panama City, FL
- * Members of the DOD Photovoltaic Review Committee

N93-25598

537-33
150507

P. 8

HIGH SPEED SOLID STATE CIRCUIT BREAKER

Thomas F. Podlesak
Electronics Engineer

U.S. Army Research Laboratory, Electronics and Power Sources
Directorate

Pulse Power Components Branch
ATTN: AMSRL-EP-MC(Podlesak)
Fort Monmouth, NJ 07703-5601

ABSTRACT

The U.S. Army Research Laboratory, Fort Monmouth, NJ, has developed and is installing two 3.3 MW high speed solid state circuit breakers at the Army's Pulse Power Center. These circuit breakers will interrupt 4160V three phase power mains in no more than 300 microseconds, two orders of magnitude faster than conventional mechanical contact type circuit breakers. These circuit breakers utilize Gate Turnoff Thyristors (GTOs) and are currently utility type devices using air cooling in an air conditioned enclosure. Future refinements include liquid cooling, either water or two phase organic coolant, and more advanced semiconductors. Each of these refinements promises a more compact, more reliable unit.

INTRODUCTION

The U.S. Army Pulse Power Center, located at Fort Monmouth, NJ, has the unique mission of performing research, development and benchmark testing of megawatt class power components and subsystems. The Center, shown in Figure 1, has an installed continuous capability of 30 MVA, continuous liquid cooling of 10 Megawatts and is shielded to the DoD TEMPEST standard. This facility is ideally suited to the stressing of high powered components to the limits of their ability.

When subjected to the extreme conditions of benchmark testing, components and subsystems are very likely to undergo a fault condition. The amount of available power makes the destruction of the component under these conditions very probable. The incident rate of component destruction, which can be hazardous to equipment and personnel, is high. In the past, the component under test was protected by standard, mechanical, utility circuit breakers, which clear in tens of milliseconds. This length of time, although considered very short under normal conditions, can be an eternity under dangerous fault conditions. It was therefore determined that a faster fault protection device was needed.



Figure 1: The U.S. Army Pulse Power Center at Fort Monmouth, NJ

The decision was made to install a solid state based circuit breaker in the 4160 Volt utility mains leading into the Pulse Power Center's dual 3.3 Megawatt power supply. This supply consists of two independent 3.3 Megawatt power supplies, which may operated together to form one 6.6 Megawatt power supply. This a most versatile unit and is used extensively for experiments at the Pulse Power Center. As part of its ongoing mission in high power components, extensive work has be done at the Pulse Power Center in the area of high power solid state switching. A series of experiments established the applicability of operating Gate Turnoff Thyristors (GTOs) in series to produce a repetitive opening and closing switch, capable of operating at voltages that were much in excess of the rating of an individual device. Additional experimental work has been conducted utilizing more advanced semiconductors, which promise more efficient operation. The knowledge of the capabilities of these devices lead to the conception of a solid state circuit breaker for power applications, with the key requirement that the interrupt time be less than 300 microseconds, two orders of magnitude better than existing mechanical circuit breaker. There is also the advantage of no mechanical degradation of the solid state circuit breaker, since moving parts have been eliminated and arcing is not present.

THE INSTALLATION OF THE SOLID STATE CIRCUIT BREAKERS

The installation of the solid state circuit breaker and associated control circuitry was authorized in 1989, under the

Productivity and Capital Improvement Program (PCIP). This is a program within the Department of Defense to improve productivity at its facilities by means of the installation of advanced equipment. The reasoning was that so much would be saved in time and material by preventing destruction of components under test that the project would easily justify the its expense, this pay back being a requirement of the program.

The specifications for the solid state circuit breaker are presented in Table I. Upon release of the specifications, the Westinghouse Science and Technology Center, Pittsburgh, PA, bid and was awarded the contract to design and manufacture two such units, one for each of the two 3.3 Megawatt power supplies. The units are illustrated in Figure 2. Note the three individual structures, one for each phase of the power circuit. The units use four GTOs per polarity per phase, resulting in twenty-four GTOs per units. The Pulse Power Center is heavily committed to preserving the technical manufacturing base in the United States and, therefore insisted on domestic GTOs. The GTOs were manufactured to Westinghouse's specifications by the Static Power Control Operation of General Electric, located at Malvern, PA. The devices are a modification of an existing GE design for 4000 V, 1 KA devices. The modification was to allow a more reliable turnon of the power, which was predicated by the control algorithm adopted by Westinghouse. The GTOs are symmetric. The devices stand approximately 4 feet (1.3 m) high, are approximately 6 feet (1.8 m) long and 2 feet (0.6 m) deep, and weigh approximately 1 ton (900 kg). They are utility units and air cooled, hence the large size and weight.

Table 1

Solid State Circuit Breaker Technical Requirements

Withstand Voltage-Terminals	15 kV
Withstand Voltage to Ground	15 kV
BIL (Basic Impulse Level)	30 kV
Continuous Current	600 A
Maximum Interrupt Current	2 kA
Surge Current (10 ms- 1 cycle)	10 kA
Steady State Switch Impedance	50 mOhms
Closing Time	50 microseconds
Opening Time	300 microseconds
Operating Temperature	0 C to 80 C

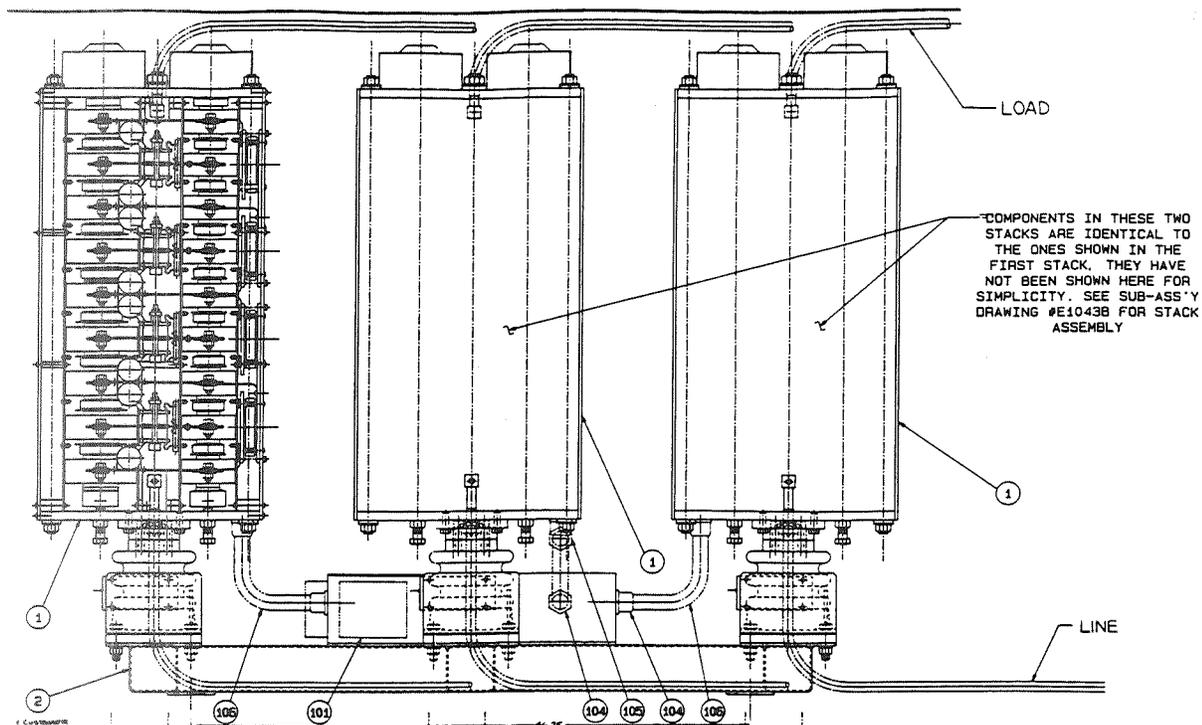


Figure 2: Drawing of the solid state circuit breaker, showing the location of the line and load cables, the control enclosure (rectangular, lower center) and detail of the structure of one phase of the circuit breaker (at left). Specifically, the eight flattened rectangles to the left of the structure are the GTOs, the corresponding eight structures to the right are steering diodes and the components along the centerline constitute the snubbers for the GTO. The unit is approximately six feet long.

The air cooling requirement necessitated a special enclosure for the units. Power handling equipment of this type is normally located outside in a switch yard. The solid state circuit breakers are no exception. They have been installed in a climatically controlled NEMA Class 4 enclosure in the switchyard located behind the Pulse Power Center, in close proximity to the existing equipment that comprise the 6.6 MW power supply. A photograph of the installed units appears as Figure 3. A schematic of the installation is shown in Figure 4. Note the voltage arrestors on both line and load sides of the solid state circuit breaker in the schematic, and which are visible at the right of the photograph. These protect the GTOs from voltage transients, which may occur in closing and opening operations. Saturable core reactors in the line side protect the unit from damage due to the instantaneous application of current. Mechanical contactors are provided on the line side for isolation. Synchronization signals from the phases are detected from the auxiliary power circuit. This is present even if the prime power is not connected. To conclude the discussion on the circuit breaker, there are three control lines, two fiber optics

and one twisted pair. The fiber optics are system enable and the on/off signal. The twisted pair is a relay closure, confirming circuit breaker closure. The signal leaves the circuit breaker over twisted pair, but is converted to fiber optic, as are all control signals on all power supplies at the Pulse Power Center.

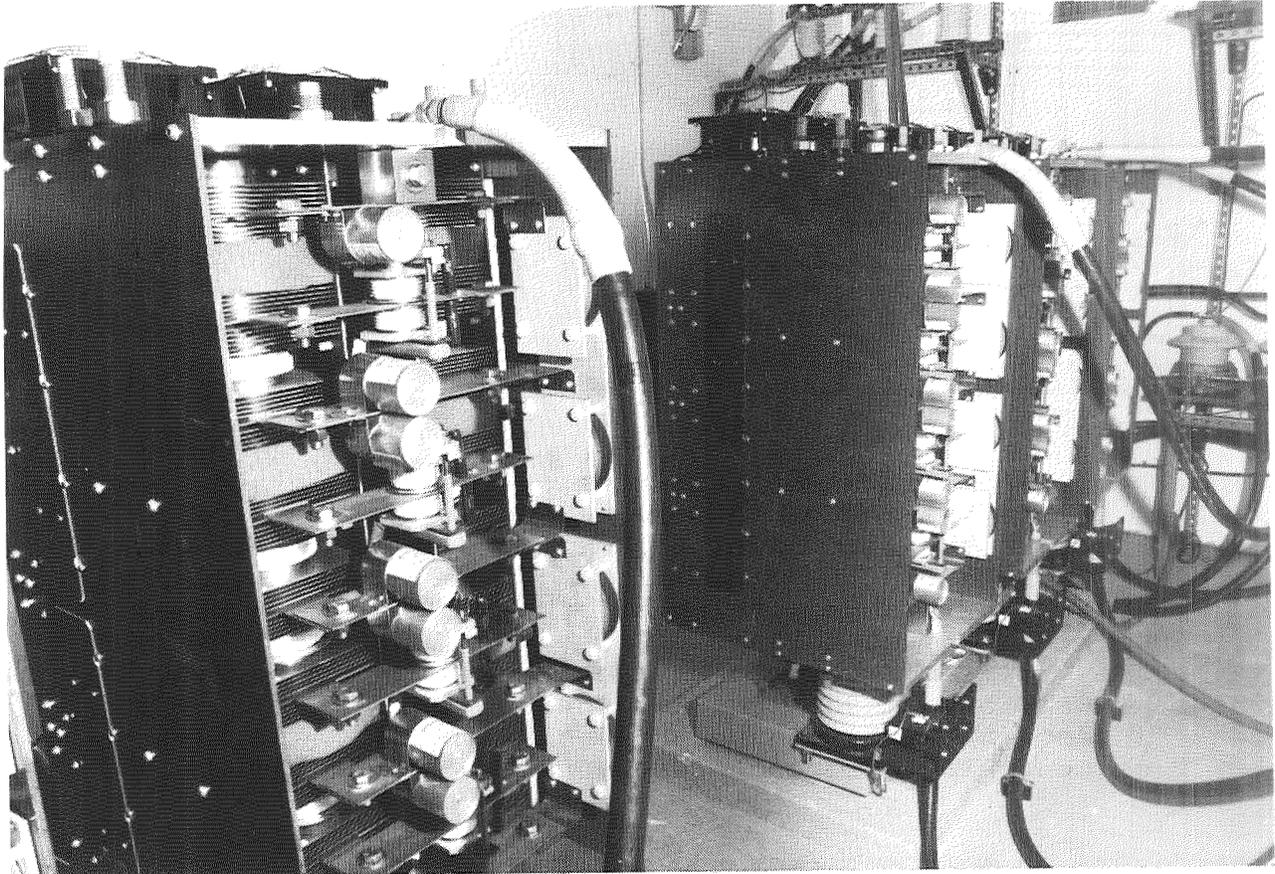


Figure 3. The solid state circuit breakers installed at Fort Monmouth. The phase section at the left clearly shows the snubber components along the vertical centerline, with the GTOs (white structures) to their left. One voltage arrester is visible to the far right.

Further on the subject of controls, the installation of the solid state circuit breakers into an existing system has provided the opportunity to upgrade the instrumentation of this system. System control, which formally consisted of mechanically driven contactors and transformer taps, has now been supplemented by a computer based control and measurement system. The computer is an Apple Macintosh IIx running the National Instruments LABVIEW instrumentation software package. The computer is interfaced to the existing control system via four plug in boards, which support various digital I/O, D to A and A to D function via new interface chassis, which have been installed in the existing control rack cabinets of the power supply. The computer also interface, via IEEE 488 bus, to two CAMAC systems, which, at the present are hosting seventy-six channels of data acquisition.

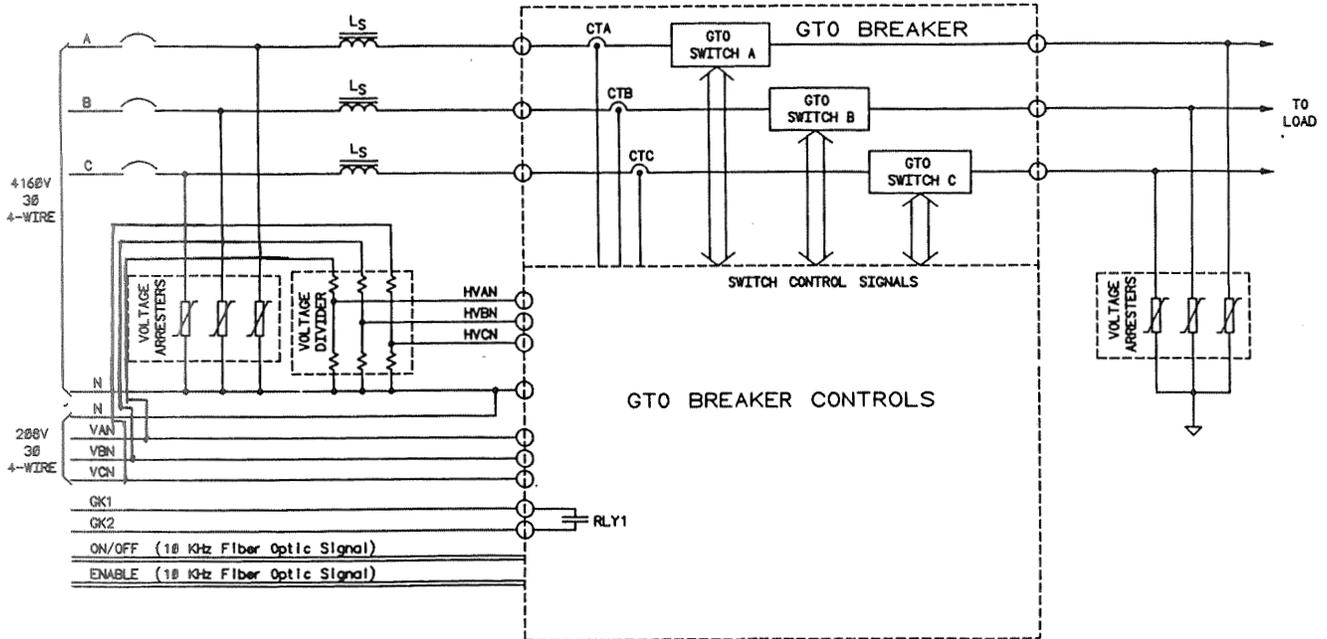


Figure 4: Schematic of installation of the solid state circuit breaker. Note the saturable core reactors and mechanical disconnects in the line side, to the left, and the voltage arrester in both line and load sides. Note also the voltage dividers, which measure phasing off the 208 V control power. Original plan had the voltage dividers measuring off the 4160 V mains, but it is possible to have no voltage on the mains, due to the aforementioned mechanical disconnects, and still need to sense phase. A thirty degree phase shift between the 4160 V and 208 V circuits is factored into the control algorithm.

The desire to add, subtract and modify this data acquisition capability lead to the selection of the CAMAC standard for this purpose. The instrument suite is completed by a custom peak/fault detector, which records peak currents and fault occurrences and will trip the system if fault number or rate per specified unit time exceeds operator set limits. The system was designed and built by Maxwell Laboratories, San Diego, CA.

STATUS OF PROGRAM

The solid state circuit breakers were completed and installed in 1991. The high voltage wiring installation was done by the Army's 535th Combat Engineering Detachment, which has a contingent based at Fort Monmouth. This military unit specializes in the installation of electric generation and distribution on the battlefield, utilizing 1 MW all fuel turbine driven 4160 V generators mounted on trailers, palletized distribution transformers and all the associated cabling. An deployed Army division requires 5 MW of electric power generation.

The new controls for the system were installed in 1992. A control problem delayed final testing of the system. However, preliminary results are available and are illustrated in Figure 5. Note the bottom trace; it is the control signal. The abrupt end of this signal is the trip command. The top trace, the current, drops to zero within 300 microseconds, exclude some small ringing out which is attributable to load conditions.

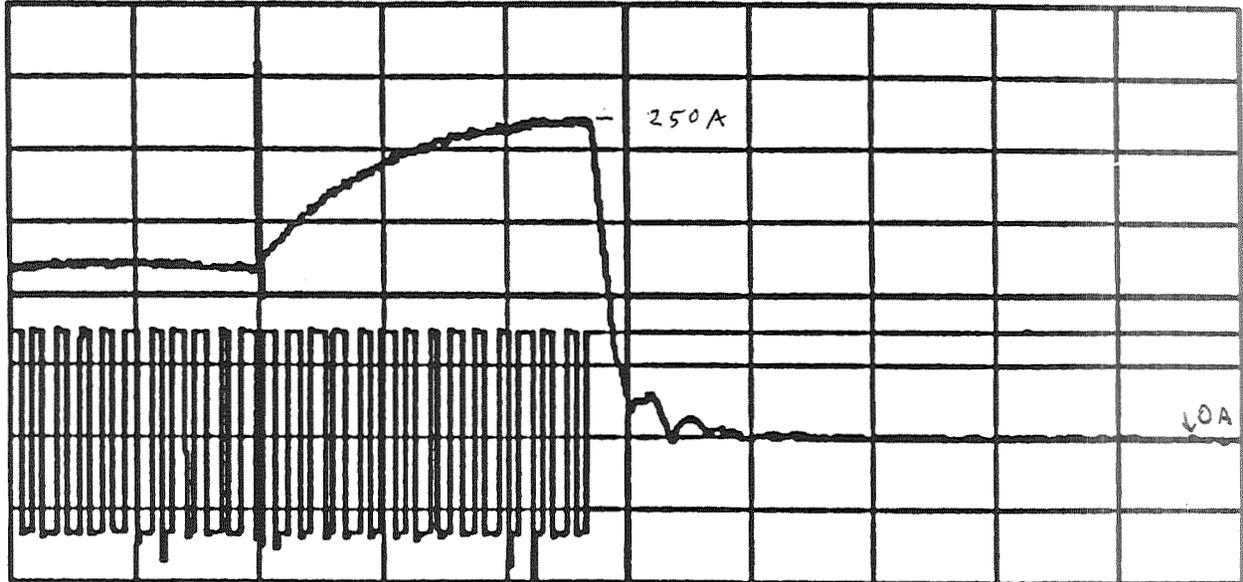


Figure 5. A test trip of one solid state circuit breaker. The top trace is the circuit breaker current and one vertical division is equal to 50 Amperes. One horizontal division is equal to 500 microseconds. The bottom trace is the 10 kHz control signal, one vertical division is equal to 5 volts. At 1 millisecond into the event, an over current is induced into the system by shorting out a portion of the load resistance. At approximately 2 milliseconds into the event, the current exceeds the trip value. The control signal ceases at 2.8 milliseconds into the event and the solid state circuit breaker trips in approximately 300 microseconds. The circuitry producing the control signal cessation is one of the slower modes, in this case; faster responses for this signal drop are possible.

FUTURE APPLICATIONS

The present application was for a utility type device, in which size and weight is of little consequence. Future units should be made much smaller. There are military applications for units of this nature. These would be in connection with the previously mentioned 1 MW trailer mounted generators. All three services use these systems. The Army, as previously mentioned, uses these for field power generation. The Navy uses them to power ships in port and the Air Force uses them to provide power to aircraft and associated equipment on the flight line. Since these operations must be portable on a worldwide basis, it is imperative that such units be lightweight, compact and rugged.

The latter requirement makes the solid state circuit breaker ideally suitable, in that it will replace the present day system which is based on mechanical contactors, often housed in fragile glass vacuum envelopes.

Civilian applications are also foreseeable for this technology, including protection of utility lines, the realization of a nonmechanical reclosure, with remote, adjustable settings, the protection of valuable industrial equipment, the improvement of power quality and balancing, improved synchronization of capacitor banks for power correction, diverters for industrial loads, power brokerage from one utility to another and intelligent power control in industrial systems. Again, these may be done at the megawatt level.

The solid state circuit breaker will achieve reduction in size and weight by several technical innovations. The first is the replacement of the current air cooling with one or two phase liquid cooling systems. Preliminary discussions with the builder of the current solid state circuit breaker indicates that size reduction, with incumbent weight reduction, to one or two cubic feet of volume for a unit rated similarly to the present units is quite conceivable. The second innovation will be a replacement of the GTOs with more advanced semiconductors. Recently developed devices have less demanding driver and snubber requirements, which will result in smaller auxiliary circuits. Further improvement may be obtained due to the more efficient nature of these devices, resulting in fewer thermal management requirements due to enhanced performance in conduction and switching.

CONCLUSION

The solid state circuit breaker is about to begin a new era in fault protection for power systems. Future more advanced systems will open up new applications both terrestrial and space based. Such devices will insure less damage due to the shortening of the time during which a fault condition exists and higher reliability due to the elimination of moving parts.

ACKNOWLEDGMENTS

The author would like to acknowledge the efforts of the two principals on this project, Dr. F. Owen Johnson of the Westinghouse Science and Technology Center and Dr. George Schofield of Maxwell Laboratories. A note of special thanks goes to Chief Warrant Officer Borst, Staff Sergeant K. Payne and Sergeants K. Chaney, C. Claxton, A. Potter, D. Brown, R. Simpson and J. Moss of the 535th Combat Engineering Detachment.

VARIABLE-SPEED GENERATORS WITH FLUX WEAKENING

150 508

A.A. Fardoun, E.F. Fuchs
Department of Electrical and
Computer Engineering
University of Colorado
Boulder, CO 80309

P.W. Carlin
National Renewable Energy Laboratory
Golden, CO 80401

P-10

ABSTRACT

A cost-competitive, permanent-magnet 20 kW generator is designed such that the following criteria are satisfied: an (over) load capability of at least 30 kW over the entire speed range of 60-120 rpm, generator weight of about 550 lbs with a maximum radial stator flux density of 0.82 T at low speed, unity power factor operation, acceptably small synchronous reactances and operation without a gear box. To justify this final design four different generator designs are investigated: the first two designs are studied to obtain a speed range from 20 to 200 rpm employing rotor field weakening, and the latter two are investigated to obtain a maximum speed range of 40 to 160 rpm based on field weakening via the stator excitation. The generator reactances and induced voltages are computed using finite element/difference solutions. Generator losses and efficiencies are presented for all four designs at rated temperature of $T_r=120^{\circ}\text{C}$.

INTRODUCTION

Most generator choices for wind turbines sized in the 50-3000 kW range are of the synchronous, squirrel-cage induction and wound-rotor induction machine type. These generators have limitations on speed range (constant, or $\pm 20\%$), and some have a relatively low power factor (e.g., p.f. < 0.9) on the machine and power system sides of the plant [1,2]. Others attempt to increase the speed range by operating a wound-rotor induction machine as a doubly fed generator [3]; the cost and complexity of this type make it undesirable. Also, current harmonics on the power system side are not controlled which degrade the quality of power, therefore, requiring passive filters and power factor correction capacitors. Moreover, all wind generators that have been built rely heavily on gear boxes. Since the gear box of a wind power plant needs to be maintained at appreciable cost about every five years [4], and is an expensive part of a wind power system, it would be advantageous to build a plant without one.

In Figure 1, a novel variable-speed wind power train is introduced. It will have the following characteristics:

A permanent-magnet machine (PMM) with a speed range of 60-120 rpm with minimum weight, size and number of poles. A newly developed buck-type rectifier [5] --which maintains a high power factor as well as a low total harmonic current distortion-- utilizing only one active switch will be used to rectify the generator output voltage. The dc voltage will then be converted to ac via an inverter. The output inverter current is controlled to produce a desired power factor as well as any desired harmonic content and thus acts as an active filter. In addition, the inverter current will be controlled such that noninteger harmonics do not exist and selected harmonics are eliminated to prevent resonance phenomena.

Since PMMs have inherent advantages over induction and synchronous machines --such as higher efficiency, no brushes, no excitation losses and no field coils or virtually no rotor losses [6]-- a permanent-magnet generator is chosen for this wind power plant application. Four different permanent-magnet machine designs generating a minimum no-load voltage of $V_{L-L}=\sqrt{3}\cdot 185\text{ V}=320.4\text{ V}$, are explored. The speed range for these machines varies from 60-120 rpm (design #3) to 20-200 rpm (design #2). The main objective of this paper is to investigate and design a PMM for wind power applications. The PMM should generate 1.5 times the rated output power over the desired speed range. To produce the maximum real output power the machine is operated at unity power factor. To approximately optimize the PMM design, the "chosen" configuration is compared with three other designs having different speed ranges and control schemes. A comparison of the four designs based upon magnet costs (volume), efficiency, weight and machine size is performed.

DESIGN OF 20 kW GENERATOR WITH FLUX WEAKENING

Two alternative three-phase permanent-magnet machine types for the generator design are considered. In the first one the flux weakening is achieved via an additional rotor coil counteracting the permanent-magnet excitation as speed increases. For this type two alternative designs (design # 1 and design # 2) are compared. The reason for studying two different designs is to investigate the possibility of a machine with a 1:10 speed range.

These designs differ in the location of the permanent-magnet as explained later. In the second type the flux weakening is controlled via field orientation of the stator excitation; no rotor coils are needed resulting in a less expensive and lighter machine. Again two alternative designs (design # 3 and design # 4) are investigated. For both designs no limits are set neither for the flux weakening nor the speed range, however, they are investigated in a way to explore the maximum possible speed range within the given limits of the terminal voltage. Designs 3 and 4 differ in the field orientation of the stator current. In design # 3 the stator current distribution is controlled to produce a unity power factor without weakening but strengthening the permanent-magnet excitation; while in design # 4 the stator current distribution is controlled to weaken the permanent-magnet excitation disregarding the unity power factor constraint. Because neodymium-iron-boron permanent-magnet material has a larger coercivity and a greater energy product than other permanent-magnet materials and thus a higher power/weight ratio [6], it is used for this generator. Since the machine speed is proportional to the phase voltage (Eq. 1), using flux weakening will bound the terminal voltage (V_t) to $V_{rated} \leq V_t \leq 2 \cdot V_{rated}$; where V_{rated} is the rated generator phase voltage at low speed.

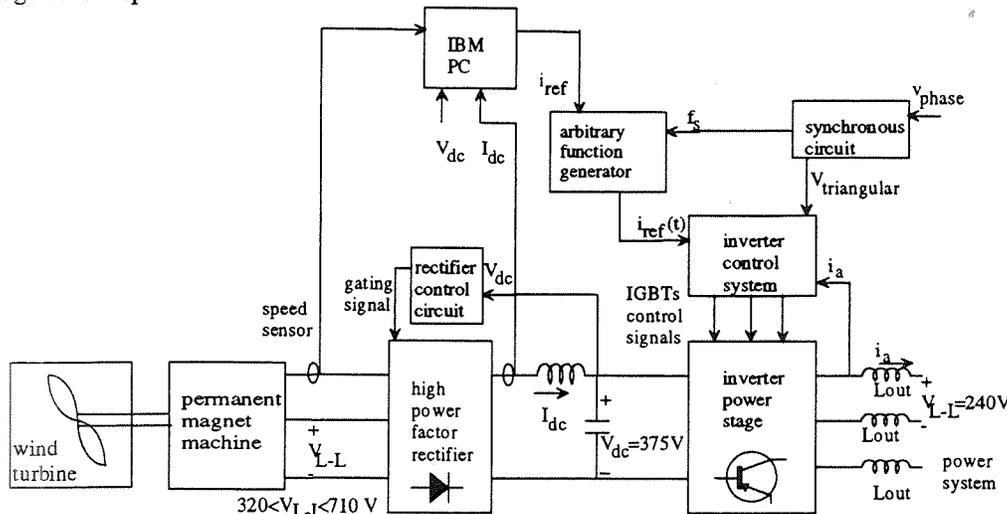


Figure 1: Variable-speed wind power plant.

$$E_{rms} = 4.44 \cdot N_{ph} \cdot f \cdot \phi_{Fe} B_{maxw} \cdot k_w \cdot l \cdot \tau_p, \quad (1)$$

where N_{ph} is the number of series turns per phase of the stator winding, f is the frequency of the machine current and voltage (e.g., $f=4$ Hz at speed of 20 rpm), ϕ_{Fe} is the iron stacking factor =0.9, B_{maxw} is the maximum (fundamental) stator radial flux density at the radial center of the stator winding, k_w is the pitch factor, l is the actual machine length, τ_p is the pole pitch length= $2\pi R_{sw}/P$, where P is the number of poles and R_{sw} is the stator radius at the radial center of the stator winding.

All designs have common constraints derived from the minimum (rated) output voltage of the machine (Figure 1).

Design Constraints

The generator design is based on the following constraints:

output power: $P_{out} = 20$ kW,

line-to-neutral voltage at no load: $V_{L-N} = 185$ V

base impedance: $Z_{base} = 5.14 \Omega$,

current density: $3.8 < J < 5$ A/mm²

copper fill factor: $0.5 < k_c < 0.75$

winding pitch: $k_w = 1.0$ (full pitch),

overload capacity: $P_{max} > 1.5 P_{rated}$ over the desired speed range,

synchronous quadrature and direct-axis reactance at low speed: $X < 1.0$ p.u.,

no gear box, forced ventilation, drip-proof,

terminal voltage limited to $V_{rated} \leq V_t \leq 2 \cdot V_{rated}$ over desired speed range,

maximum rated temperature: $T_{rated} = 140^\circ\text{C}$.

Generator Design: Type # 1

For this type two alternative designs are considered. The first one, design # 1, has permanent magnets mounted on the surface of the rotor pole shoes and the flux weakening coil is located --like in the case of a

conventional synchronous machine-- in the interpolar rotor space (Figure 2a). The second design, design # 2, has permanent magnets mounted within the rotor yoke and the flux weakening coil is located in the interpolar space as for design # 1 (Figure 2b). However, the flux weakening coil lies between the stator winding and the permanent-magnet excitation for design # 2 while for design # 1 the permanent magnet lies between the stator winding and the flux weakening coil.

From a linear analysis point of view -neglecting saturation and leakage effects- both designs are the same. However it has been shown in [7] that design # 2 provides a wider speed range than design # 1 due to the nonlinear behavior of the machine caused by saturation and leakage. Since the linear analysis neglects iron-core saturation and the leakage flux of the machine, its validity is severely limited. The linear approach will reveal linear flux weakening characteristics, while the nonlinear analysis, based on numerical solutions (e.g., finite element/ difference techniques), indicates that the characteristics are nonlinear.

Numerical magnetic field solutions for the generator cross-section are obtained with an available finite element/difference software program. This software is modified to permit the modeling of permanent-magnet machines. The magnet is approximated by thin current sheets along both sides of the permanent magnet [8]; the magnet behaves almost like an air gap with a permeability of $\mu_m = 1.06\mu_0$. The flux weakening coils (ampere-turns) are assumed to be uniformly distributed within the rotor slot and modeled by inputting the desired ampere-turns in the rotor slot (Figure 2a). In addition, inputs to the field calculation program are the magnetic characteristics (e.g., iron-core characteristics, magnet height, coercivity and flux weakening current) and the geometric dimensions of the machine. The output consists of the vector potentials and the radial flux density of the machine within one pole pitch. Given the flux density in the radial center of the stator, the PMM phase voltage is computed using Eq. 1.

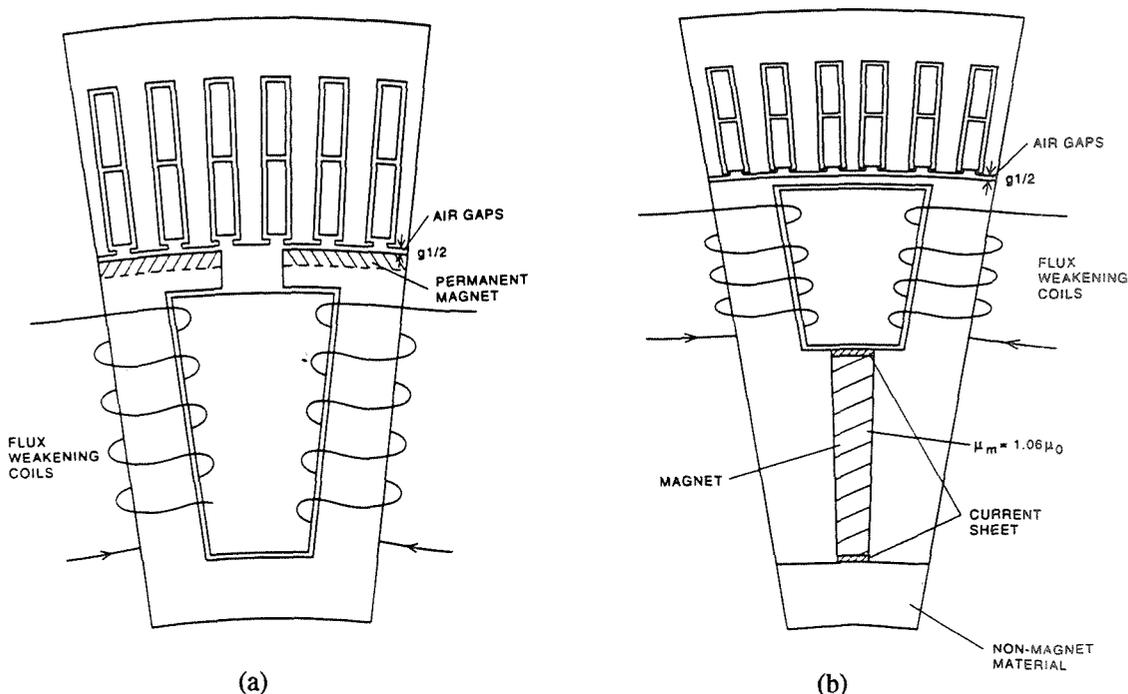


Figure 2: Geometric dimensions of designs # 1 (a) and # 2 (b).

Design # 1: For design # 1 the permanent magnets of the machine are mounted on the surface of the rotor pole shoes as shown in Figure 2a. The flux weakening characteristics at no load and full load are computed via finite element/difference method. Figures 3a, b illustrate flux distributions of design # 1 at full load for flux weakening of -500 and -6700 AT, respectively. Flux weakening at no load is possible down to about $B_{maxw} = 0.3$ T and at full load down to 0.45 T. The flux weakening current cannot be increased because the permanent magnet might be demagnetized. Figure 4 presents the flux weakening characteristics of the PMM versus speed and B_{maxw} for design # 1 at full load. In this figure B_{maxw} is the maximum radial flux density within the radial center of the stator winding and $(AT)_{FW}$ are the total ampere turns of the flux weakening coil within one rotor slot. Since the flux weakening coil is not situated in between the stator winding and the permanent magnet,

armature reaction and leakage flux effects limit the flux density reduction to 0.45 T and one achieves a speed range of 1:4.

Design # 2: In design # 2 the permanent magnet is placed within the rotor yoke as shown in Figures 2b. In this configuration the flux weakening coil lies in between the stator winding and the permanent magnet. Nonlinear numerical solutions show that this location of the magnet makes it possible for the flux weakening coil to reduce the radial flux density in the stator to virtually zero at no load condition and to 0.2 T at full load. Figures 3c, d show flux distributions for design # 2 at full load for speeds of 20 and 200 rpm, respectively. Computed characteristics at full load are given for design # 2 in Figure 4. It should be noted that as flux weakening increases the value of the quadrature reactance (X_Q) increases as explained in [7].

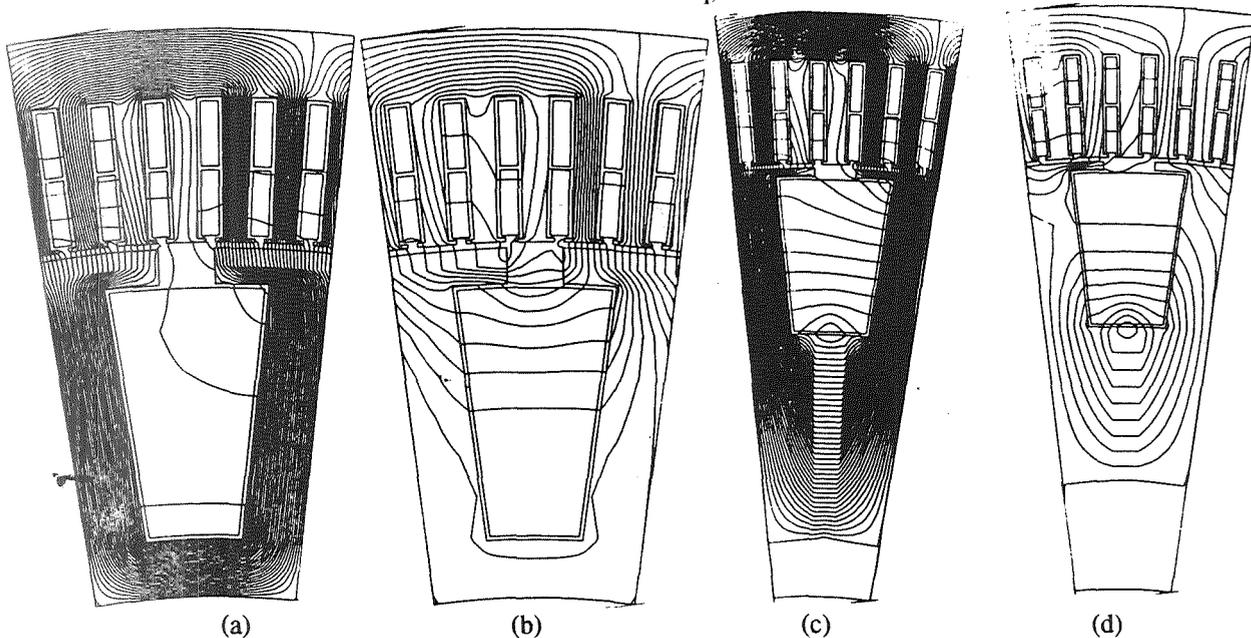


Figure 3: Field distribution for design # 1 at full load for (a) $n=20$ rpm at $(AT)_{FW}=-500$ AT, (b) $n=80$ rpm at $(AT)_{FW}=-6700$ AT, and design # 2 for (c) $n=20$ rpm at $(AT)_{FW}=-500$ AT, (b) $n=200$ rpm at $(AT)_{FW}=-6600$ AT .

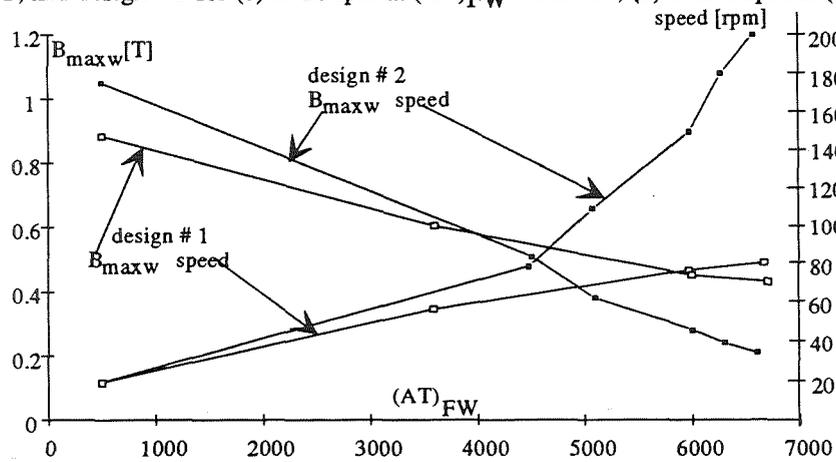


Figure 4: $(AT)_{FW}$ versus B_{maxw} and speed at full load (designs # 1 and 2).

Generator Design: Type # 2

In this type flux weakening is controlled by varying the phase angle of the stator magnetomotive force (mmf) with respect to the rotor excitation. Two alternative designs (designs # 3 and 4) are investigated. The magnet excitation is weakened by controlling the phase angle of the stator mmf and the amplitude of the stator current such that the output power is equal to or greater than 20 kW.

Design # 3: In design # 3 the stator current is controlled such that it is in phase with the stator voltage (unity power factor operation). The permanent-magnet excitation is not opposed, on the contrary the flux density in the air gap increases due to the stator current distribution which reinforces the permanent-magnet mmf as

shown by the phasor diagram. Because the voltage increases due to the strengthening of the mutual flux by the stator mmf, the machine current is below its rated value. This results in a reduction of the stator slot size and therefore weight. Since the flux density is not decreased (no flux weakening), the speed is limited to a 1:2 range because the terminal voltage is limited by the same ratio to $V_{\text{rated}} \leq V_t \leq 2 \cdot V_{\text{rated}}$. To maintain a constant output power, the machine current is decreased at high speed to half the current value at low speed since the voltage increases by a factor of 2 at high speed. Hence, the losses at high speed are less than those at low speed. It is also important to note that the weight of this machine is significantly smaller than those of type # 1 due to the fact that no flux weakening is applied and no rotor coils are used. Since this design is less expensive --as the number of poles decreases, it will be easier to manufacture and the magnet volume decreases-- it will be investigated for 12, 16, 20, and 24 poles. However, the number of stator turns increases in order to guarantee the minimum no-load voltage. The cost, weight, efficiency, magnet volume and machine reactances for 12, 16 and 20 poles are presented in Table 1.

Full-Load Analysis: The performance of design # 3 at full load is investigated using a finite element/difference software. The machine load is modeled by the stator current and the angle ψ [7], which is the angle between the stator mmf F and the fictitious induced voltage \tilde{E}_o . Since one of the constraints for this design is to operate at unity power factor load ($\theta=0$) ψ is equal to the torque angle δ . The phasor diagram at $n=60$ rpm is presented in Figure 5. Note that the quadrature machine reactance X_q can be computed using the phasor diagram. Both the reactances determined from flux linkages and those obtained from the phasor diagram match very well.

$$X_q = (V_t / I_t) \tan \delta \quad (2)$$

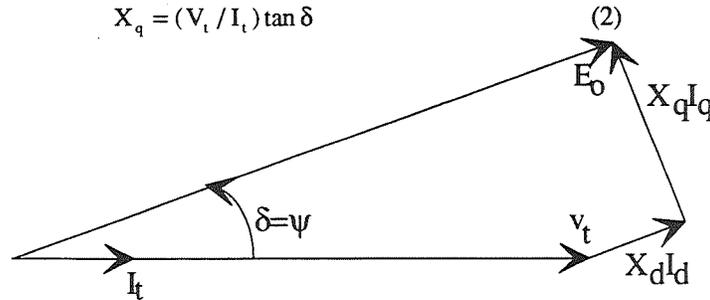


Figure 5: Phasor diagram for speed of $n=60$ rpm.

Flux distributions for design # 3 at no load and full load at $n=60$ rpm are shown in Figure 6a and 6b, respectively. As mentioned before, due to the stator current excitation which strengthens the permanent-magnet mmf resulting in a higher radial flux density at unity power factor load than at no load. The strengthening of the mutual flux increases the flux density B_{maxw} to about 0.82 T (from 0.63 T at no load). This is fortuitous since the utilization of the machine increases with increasing stator load currents; or in other words the no-load induced voltage is always somewhat smaller than the induced voltage at any load at unity power factor.

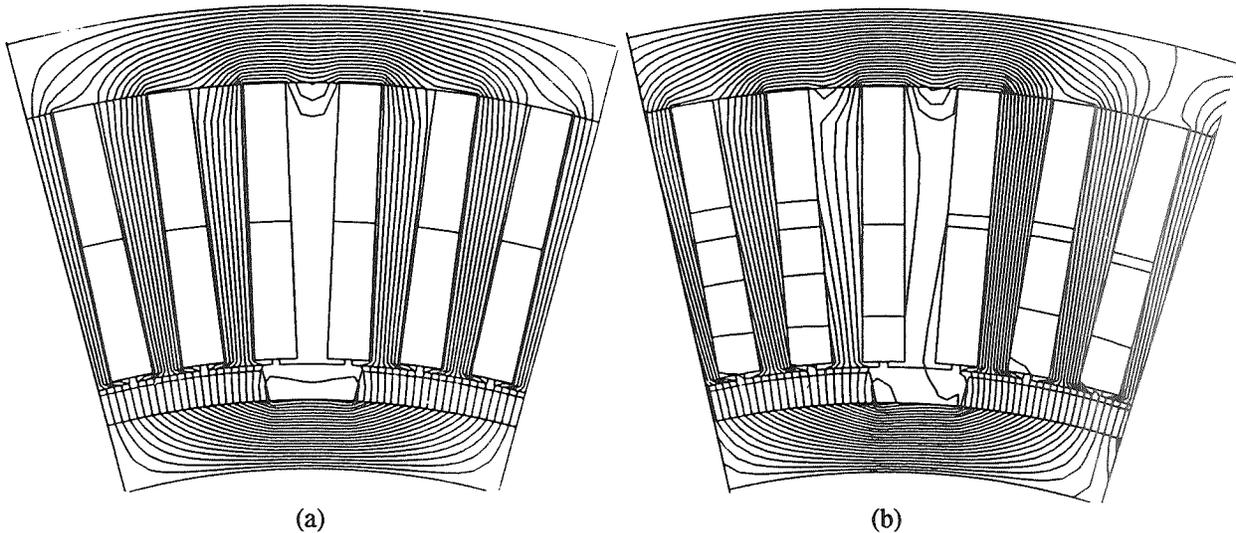


Figure 6: Flux distribution for design # 3 (a) at no load and (b) full load $n=60$ rpm.

Calculation of Direct and Quadrature Synchronous Inductances: The direct (L_d) and quadrature (L_q) inductances of one phase of the stator winding of the machine are defined as

$$L_d = L_{md} + L_{dl} + L_{de} \text{ and } L_q = L_{mq} + L_{ql} + L_{qe}, \quad (3a, b)$$

where L_{md} , L_{dl} and L_{de} are the mutual, leakage and end leakage direct inductances, respectively; L_{mq} , L_{ql} and L_{qe} are the mutual, leakage and end leakage quadrature inductances. Note that all inductance values for the direct and quadrature inductances are computed using the same flux linkage method explained below, except that $\psi=90^\circ$ for L_d while $\psi=0^\circ$ for L_q .

Synchronous Self Inductance without L_e : The synchronous inductance is computed using the flux-linkage method [9]

$$L = \Psi / I, \quad (4)$$

where Ψ are the flux linkages per phase and I is the current in one turn. The flux linkages are approximated as the sum of the flux linkages of the top and bottom layer coil sides of a phase belt. The flux linkages of a top (t) coil side with N_c turns (see Figure 7) are

$$\Psi^t = N_c \cdot A_{sm\tau}^t \cdot I \cdot \varphi_{Fe}, \quad (5)$$

where $A_{sm\tau}^t$ is the vector potential at the center of the $m\tau^{\text{th}}$ coil side of the top layer of the stator winding. The

vector potentials $A_{sm\tau}^t$ are computed from incremental magnetic stator field distributions existing at a given load: the reluctivities obtained for a given operating point (e.g., full load) are maintained constant that is "frozen" and are used for the calculation of the synchronous inductance of one stator phase. The inductances of the top and bottom coil sides of phase B of the stator winding can be approximated as defined in Figure 7 and described by Eqs. 6a and 6b, respectively.

$$L^t = \frac{1}{I_B^t} \sum_{m\tau=1}^{M\tau} N_c \cdot I \cdot \varphi_{Fe} \cdot A_{sm\tau B}^t, \quad (6a)$$

and

$$L^b = \frac{1}{I_B^b} \sum_{m\tau=1}^{M\tau} N_c \cdot I \cdot \varphi_{Fe} \cdot A_{sm\tau B}^b, \quad (6b)$$

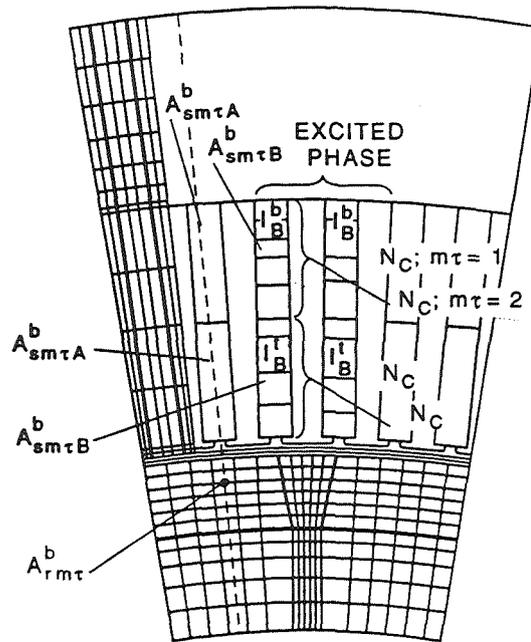


Figure 7: Definition of the stator self inductance for $M\tau=2$.

where $M\tau$ is the number of coils per pole, $A_{sm\tau B}^b$ is the vector potential at the center of the $m\tau^{\text{th}}$ bottom layer coil side of the stator winding of phase B and I_B^1 is the current in one turn of the coil side. The sum of Eqs. (6a) and (6b) is the inductance of a phasebelt L^i for $i=A, B$ and C . The direct (X_d) or quadrature (X_q) reactance of a phase can now be computed as

$$x_s = 2\pi f \cdot \beta(L^A + L^B + L^C) / \alpha \cdot m_s, \quad (7)$$

where β is the number of phasebelts in series which is equal to the number of poles P for the connection at hand, α is the number of phasebelts in parallel which is 1 for the given application, and m_s is the number of stator phases. The direct and quadrature synchronous reactance values per phase of the machine are listed in Table 1.

Armature Leakage Reactances: The leakage flux linkages are computed using Eq. 4. Since the flux linkages for the leakage inductance are those that link the stator and do not link the rotor (see Figure 7), Equation 5 is modified as

$$\Psi_l = N_c \cdot I \cdot \varphi_{Fe} (A_{sm\tau}^t + A_{sm\tau}^b - 2A_{r_{m\tau}}), \quad (8)$$

where Ψ_l are the total flux linkages (top and bottom coils) for the leakage inductance, $A_{sm\tau}^t$ and $A_{sm\tau}^b$ are the vector potentials within the top and bottom layers of the $m\tau^{\text{th}}$ coil in the stator and $A_{r_{m\tau}}$ is the corresponding rotor vector potential as shown in Figure 7. Again $\psi=90^\circ$ for L_{dl} while $\psi=0^\circ$ for L_{ql} . Assuming that the top and bottom layers of the stator coil (phase B) have the same current the leakage inductance corresponding to phase B is

$$L_l = \frac{1}{I_B} \sum_{m\tau=1}^{M\tau} N_c \cdot I \cdot \varphi_{Fe} (A_{sm\tau B}^b + A_{sm\tau B}^t - 2A_{r_{m\tau}}). \quad (9)$$

The leakage reactance is computed using Eq. 7. Typical leakage reactances at low speed for design # 3 are listed in Table 1. Figures 8a, and 8b represent field distributions for the calculation of direct and quadrature synchronous inductances as well as the armature leakage inductance.

Armature End Leakage Reactance: The end leakage reactances are approximated as

$$L_{de} = L_{dl}/3 \quad (10a)$$

and

$$l_{qe} = L_{ql}/3. \quad (10b)$$

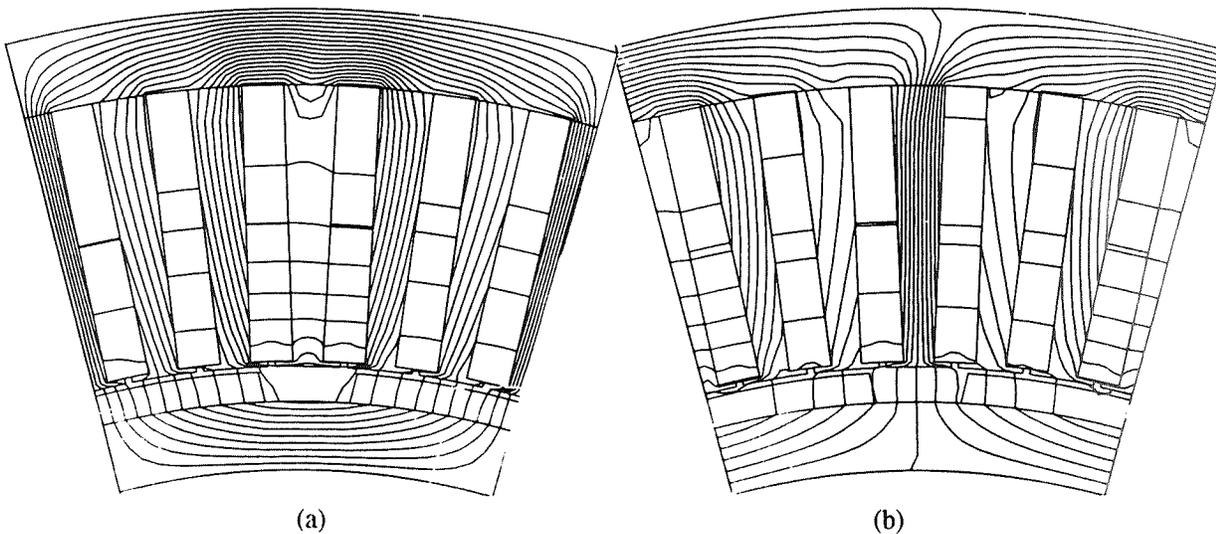


Figure 8: Field distribution for the calculation of (a) the saturated direct synchronous inductance L_d and armature leakage inductance L_{dl} (design # 3, 12 poles) and (b) saturated quadrature synchronous inductance L_q and armature leakage inductance L_{ql} (design # 3, 12 poles).

Power Calculation: The output power of the generator can be determined in two alternative ways : first from

$$P_{out} = 3V_t I_t \cos\theta, \quad (11)$$

and

$$P_{out} = 3\left(\frac{E_o V_t}{X_d} \sin\delta + \frac{X_d - X_q}{2X_d X_q} V_t^2 \sin 2\delta\right). \quad (12)$$

The results of both approaches corroborate quite well. The output power versus the torque angle δ characteristics are shown in Figure 9 for speeds of 60 and 120 rpm.

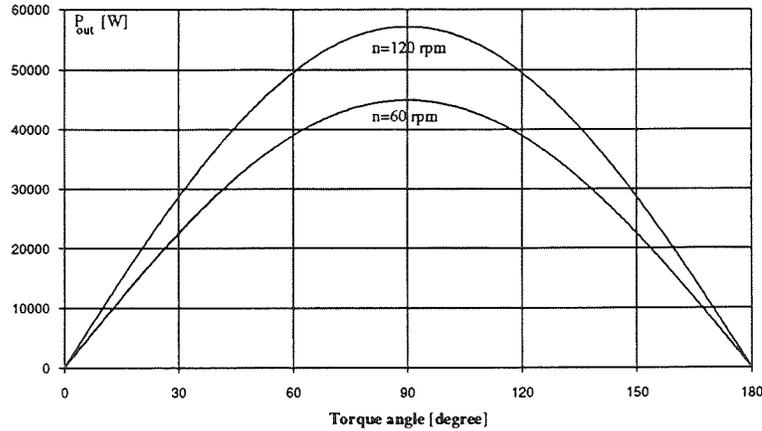


Figure 9: Power versus torque angle characteristics at low and high speeds (design # 3).

Table 1: Comparison of cost for design # 3 with 12, 16 and 20 poles.

# of poles	12		16		20	
	60	120	60	120	60	120
speed [rpm]	60	120	60	120	60	120
magnet volume per pole [cm ³]	112.2		109.8		90.9	
weight [lbs.]	537		570.7		511	
# of turns per phase per pole	32		18		14	
total losses [W]	1846	730.3	2939.9	1139.3	2637.4	1096.5
ohmic losses [W]	1795	586.2	2856.6	903.5	2533.1	801.5
core losses [W]	51	144.1	83.3	235.8	104.3	295
efficiency [%]	90.8	96.3	85.3	94.3	86.8	94.5
radial flux density [T]	0.818	0.709	0.776	0.713	0.764	0.714
L-N voltage [V]	235.9	408.9	207.5	381	203.8	379
phase current [A]	28	16	32	18	32	18
d-axis reactance [Ω]	3.36	6.71	2.1	4.2	1.55	3.1
d-axis leakage reactance [Ω]	1.16	2.32	0.7	1.4	0.6	1.2
q-axis reactance [Ω]	2.51	3.02	1.7	3.4	1.45	2.9
q-axis leakage reactance [Ω]	1.09	2.18	0.68	1.36	0.55	1.1

Design # 4: In this design the permanent-magnet excitation is weakened by controlling the relative position of the stator mmf with respect to the rotor mmf so that the stator mmf opposes that of the permanent magnet as shown in the phasor diagram of Figure 10. Controlling the air gap flux via the stator excitation requires a six-switch rectifier instead of one switch as proposed for designs # 1, 2 and 3. Even though this machine has no flux weakening rotor coil losses, its stator conduction losses increase because the stator current is increased as the speed increases to weaken the flux. It also should be noted that the weight of this machine is less than those of machines of designs # 1 and # 2. However, the speed range of this configuration is limited to 1:4 while that of design # 2 is from 1:10.

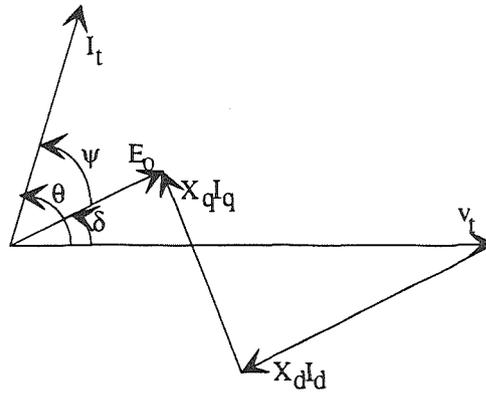


Figure 10: Phasor diagram for design # 4 at $n=160$ rpm.

DESIGN CONSIDERATION FOR 300 kW GENERATOR

The weight of the 20 kW machine is approximated as

$$W_{20\text{kW}} = \varphi_e \pi (R_{so}^2 - R_{ri}^2) l \gamma_{Fe}, \quad (13)$$

where φ_e is a filling factor and γ_{Fe} is the iron specific weight. Model laws for electrical machines indicate that the output power P_{out} of a machine relates to the dimension X as $P_{\text{out}} \propto (X)^4$ while the volume (weight) of a machine increases according to $W \propto (X)^3$. Since the present 20 kW is somewhat oversized one can extract at least 30 kW from the present volume (or weight) and therefore a 300 kW machine will weigh

$$W_{300\text{kW}} = W_{20\text{kW}} \cdot 10^{3/4} \quad (14)$$

The approximate weights for all four machine designs is listed in Table 2. However, it should be noted that it is difficult to mount the rotor of a permanent-magnet machine of 300 kW (designs 1, 3 and 4). An alternative solution is to locate the magnet within the rotor yoke as employed for design # 2, where the permanent magnets can be installed after the rotor has been mounted.

DISCUSSION AND CONCLUSIONS

Four machine designs are investigated; the first one is not recommended because design # 2 has a wider speed range with almost the same weight and efficiency. Therefore, the task remains to choose one machine from the remaining three designs for construction and testing. Table 2 lists magnet costs (magnet size), machine weights as well as efficiencies and losses. One should also consider that this application deals with wind turbines as prime movers. In practice about 80 % of the available wind energy is extracted by wind turbines with constant-speed generators [4]. With a machine of a speed range of 1 to infinity one could extract the remaining 20 %. However, a machine with a speed range of 1 to 2 (e.g., 60-120 rpm) one can extract about 18% of the available 20% [4]. This fact leads to the simple conclusion that design # 3 is the most cost-competitive design for wind power plants. Note that design # 3 with 12 poles is chosen for construction since it has 72 stator slots and the appropriate laminated core is an off-the-shelf item; in other words it is the least expensive among the 12, 16 and 20 pole configurations. The efficiency of the 12 pole arrangement is the highest among all other designs. Efficiencies for all four designs are shown in Figures 11. Note that the efficiency for design # 3 increases as the speed increases because the current decreases and the voltage increases maintaining constant output power without applying any flux weakening. Efficiencies for designs 1 and 2 are maximum around $n=40$ rpm since the terminal voltage is assumed to increase from 20 rpm to 40 rpm and the current is allowed to decrease and thus the stator losses decrease. At high speeds the losses increase due to the high rotor losses even though the stator losses decrease due to lower stator current. For design # 4, the stator losses increase at high speed since the stator current is increased to weaken the machine radial flux density. The temperature in all generators is assumed to be 120 °C and the conduction losses are increased accordingly. The magnet size for designs # 1 and 2 is adjusted assuming that the magnet temperature is 140 °C while for designs # 3 and 4 it is sized such that the magnet temperature is 120 °C. It should be mentioned that the performance of design # 4 for different pole numbers has not been researched, however, this will not change the validity of the above conclusion for choosing design # 3 as the most cost-competitive design for this wind power plant application.

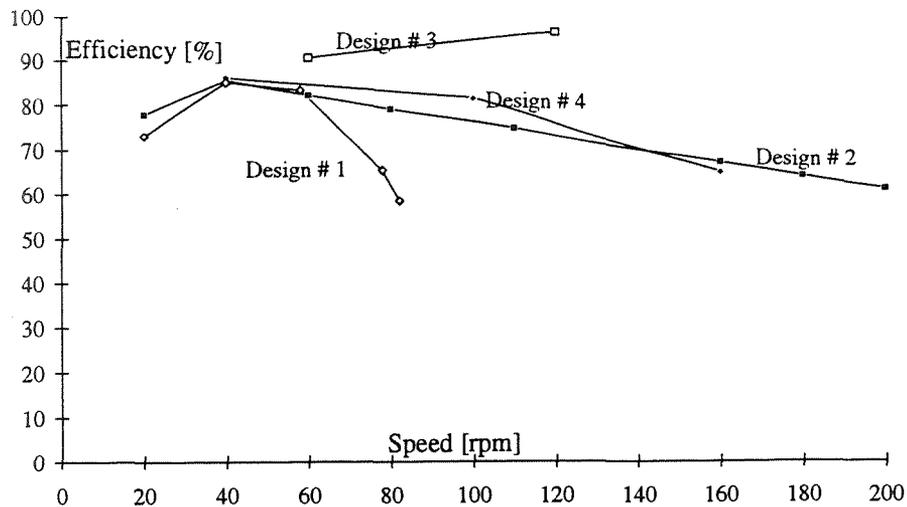


Figure 11: Efficiencies for all four generator designs.

Table 2: Comparison of cost, weight and efficiency of all four designs.

	Design # 1	Design # 2	Design # 3	Design # 4
# of poles	24	24	12	24
magnet volume per pole [cm ³]	92.	231.1	112.2	67.
weight [lbs]	1286	1584.	537	638.3
weight of 300 kW machine [lbs.]	7231.7	8907.5	3019.8	3589.4
speed range [rpm]	20-85	20-200	60-120	40-160
high speed efficiency [%]	58.6	61	96.3	64.3
total high speed losses [W]	8288	7796.2	730.3	7059.8
high speed stator losses [W]	1218	1013.1	586.2	6504.8
high speed rotor losses [W]	7070	6783.1	144.1	555
low speed efficiency [%]	73	77.7	90.8	86.5
total low speed losses [W]	5408.2	4452	1846	2703.6
low speed stator losses [W]	5366.2	4410	1795	2634.3
low speed rotor losses [W]	42	42	51	69.3
energy flow	one direction	one direction	one direction	both directions

REFERENCES

- [1] T. H. Lauw, "AC-DC-AC Conversion Systems for Mains-Connected Windpower Generation", *ASME Wind Energy Symposium*, 1988.
- [2] E. Hinrichsen, "Variable Rotor Speed for Wind Turbines, Objectives and Issues", *Windpower*, San Francisco, CA, 1985.
- [3] P. Holmes and N. Elsonbatty, "Cycloconverter-Excited Divided-Winding Doubly-Fed Machine as a Wind-Power Converter", *IEE Proceedings-B*, Vol. 131, March 1984, pp. 61-68.
- [4] C. P. Butterfield, W. Musial and P. W. Carlin, National Renewable Energy Laboratory, private communication, June 1992.
- [5] E. Ismail and R. W. Erickson, "A Single Transistor Three Phase Resonant Switch for High Quality Rectification", *IEEE Power Electronic Specialist Conference*, Spain, July 1-7, 1992, pp. 1341-1351.
- [6] P. Mellor, B. Chaaban and K. Binns, "Rare Earth Permanent-Magnet Motors Avoiding Measurement of Load Angle", *IEE Proceedings-B*, Vol. 138, No. 6, 1991, pp. 322-330.
- [7] E. F. Fuchs, A. A. Fardoun, P. Carlin and R. W. Erickson, "Permanent Magnet Machines with Large Speed Variations", *American Wind Energy Association Conference*, Seattle, 19-23 October, 1992.
- [8] J. R. Brauer, L. A. Larkin and V. D. Overbye, "Finite Element Modeling of Permanent Magnet Devices", *Journal of Applied Physics*, Vol. 55, pp. 2183-2185, March, 1984.
- [9] E. F. Fuchs, M. Poloujadoff and G. W. Neal, "Starting Performance of Saturable Three-Phase Induction Motors", *IEEE Trans. on Energy Conversion*, Vol. 3, No. 3, September 1988, pp. 624-633.



OMI

**INFORMATION AND COMMUNICATIONS PART 2:
COMPUTER SIMULATION AND MODELING**

PRECEDING PAGE BLANK NOT FILMED

364. C-5

omit

INDUSTRIAL APPLICATIONS OF COMPUTATIONAL FLUID DYNAMICS

This paper was withdrawn from presentation.

539-6/
150509
p 10

SCIENTIFIC VISUALIZATION USING THE
FLOW ANALYSIS SOFTWARE TOOLKIT (FAST)

N 93-25600

Gordon V. Bancroft
Paul G. Kelaita
R. Kevin McCabe
Fergus J. Merritt
Todd C. Plessel
Timothy A. Sandstrom
John T. West

Sterling Software Inc.
MS 258-2, NASA Ames Research Center
Moffett Field, CA 94035

ABSTRACT

Over the past few years the Flow Analysis Software Toolkit (FAST) has matured into a useful tool for visualizing and analyzing scientific data on high-performance graphics workstations. Originally designed for visualizing the results of fluid dynamics research, FAST has demonstrated its flexibility by being used in several other areas of scientific research. These research areas include earth and space sciences, acid rain and ozone modelling, and automotive design, just to name a few. This paper describes the current status of FAST, including the basic concepts, architecture, existing functionality and features, and some of the known applications for which FAST is being used. A few of the applications, by both NASA and non-NASA agencies, are outlined in more detail. Described in the outlines are the goals of each visualization project, the techniques or 'tricks' used to produce the desired results, and custom modifications to FAST, if any, done to further enhance the analysis. Some of the future directions for FAST are also described.

INTRODUCTION

With supercomputer technology advancing at a rapid pace, the amount and size of data being generated by computational researchers is becoming enormous. Tools are desperately needed to help the researchers process this data - to reduce it into manageable, understandable and useful pieces of information. For the past few years the Flow Analysis Software Toolkit (FAST) has been used by scientific researchers at NASA Ames Research Center to visually analyze their computational data. The use of FAST is becoming more widespread due to its generalized visualization approach. Originally developed for the analysis of computational fluid dynamics (CFD) results, FAST is increasingly being used by scientists to visualize data from other disciplines at research institutes across the country.

BACKGROUND

The NAS Project

In the early 1980's the Numerical Aerodynamics Simulation (NAS) project was started at NASA Ames. NAS's charter was and is to provide CFD researchers with the fastest supercomputers, the fastest networks and best support systems, including graphics workstations. On the initial NAS systems, which consisted of a CRAY 2 supercomputer, hyperchannel networks and Silicon Graphics IRIS 2000 series workstations, researchers typically generated simple 2- and 3-D volumetric data depicting fluid flow around and within solid bodies. These data sets were typically analyzed (or 'visualized') using several different software packages. In fact, the entire CFD research cycle, which includes generating the computational volume (grid), computing fluid flow characteristics within the grid, and analyzing the results, consisted of running several independent programs on both the supercomputer and workstation and moving the data files back and forth via network file transfers. This mode of operation was generally acceptable since the amount of data was fairly manageable.

The late 80's brought a new wave of supercomputers, networks and workstations that dramatically increased the amount and complexity of data being generated by researchers. It was at that time that the members of the Workstation Applications Office (WAO, part of the Fluid Dynamics division at Ames) recognized the need for a software package that could assist the researchers in all phases of the research cycle. A package was needed that was flexible, so new functionality could be added easily, consistent, so that users could learn the basics and apply that knowledge in using the rest of the package, and distributed, so that users could access and manipulate their data without having to transfer files over network lines.

Initial FAST Proposal

The WAO team proposed FAST as an 'environment' that would provide researchers with a full range of visual analysis tools, animation production capability, and the potential to monitor and/or direct their flow codes and grid generators on the fly. The target workstation was the Silicon Graphics IRIS 320 VGX - a high-performance graphics system with two processors. The original FAST design concepts were as follows:

Minimal Data Handling - Alleviate the need for the researcher to move files around.

Flexible Architecture - Implement modular design to facilitate easy addition of functionality.

Consistent User Interface - Keep the same user interface across all applications.

Highly Interactive - Provide instant feedback for each user action.

Processes Distributed to Optimal Resources - Utilize all available computing power.

The first versions of FAST achieved all of the original design goals except the distribution of processes to optimal computing resources. This capability is a planned addition to the next version of the software.

Current Status

Initial versions of FAST were released and distributed as a part of the beta test program by WAO. Currently FAST is being used at over 400 sites nationally. Development of the code is continuing under the NAS division. The latest version, which is available through COSMIC, runs on all Silicon Graphics workstations. The distribution package includes source code and a user's manual.

CURRENT ARCHITECTURE

The architecture of FAST is based on a modular design with shared memory used for storing common data. The collection of modules, each a separate UNIX process that provides a specific and independent function, makes up the FAST environment. The central module of the system is the Hub. The Hub allocates and maintains the 'pool' of data in shared memory and links all other modules to the system with two-way sockets (inter-process communication utilities). This 'hub-and-spoke' design facilitates inter-module communication as well as centralized control of the system. (See figure 1)

The remaining modules in the FAST environment are grouped into three categories: data manipulation, object manipulation and viewing. The data manipulation modules are tasked with transferring data, which typically consist of grids and flow code solutions, from the disk (or in future versions from a remote supercomputer) to core memory of the workstation. They are also responsible for computing scalar and vector quantities based on the grid and solution data. The object manipulation modules use the fundamental data to generate graphical objects such as cutting planes, isosurfaces and particle traces. The collection of these objects are then grouped, transformed and rendered to the workstation display by the viewing module which currently is the Viewer.

As stated above, this architecture fulfills all but one of the original design goals of FAST. The use of shared memory for data storage gives the user easy access to the data without having to move it from process to process. The modular design allows the system to be easily extended to include the user's own customized application. It also facilitates immediate feedback in that if the user changes a module's object the scene containing that object is

immediately updated by the Viewer. The user interface looks and feels the same in all modules so that once the user becomes familiar with basic UI concepts, learning additional components should be expedited.

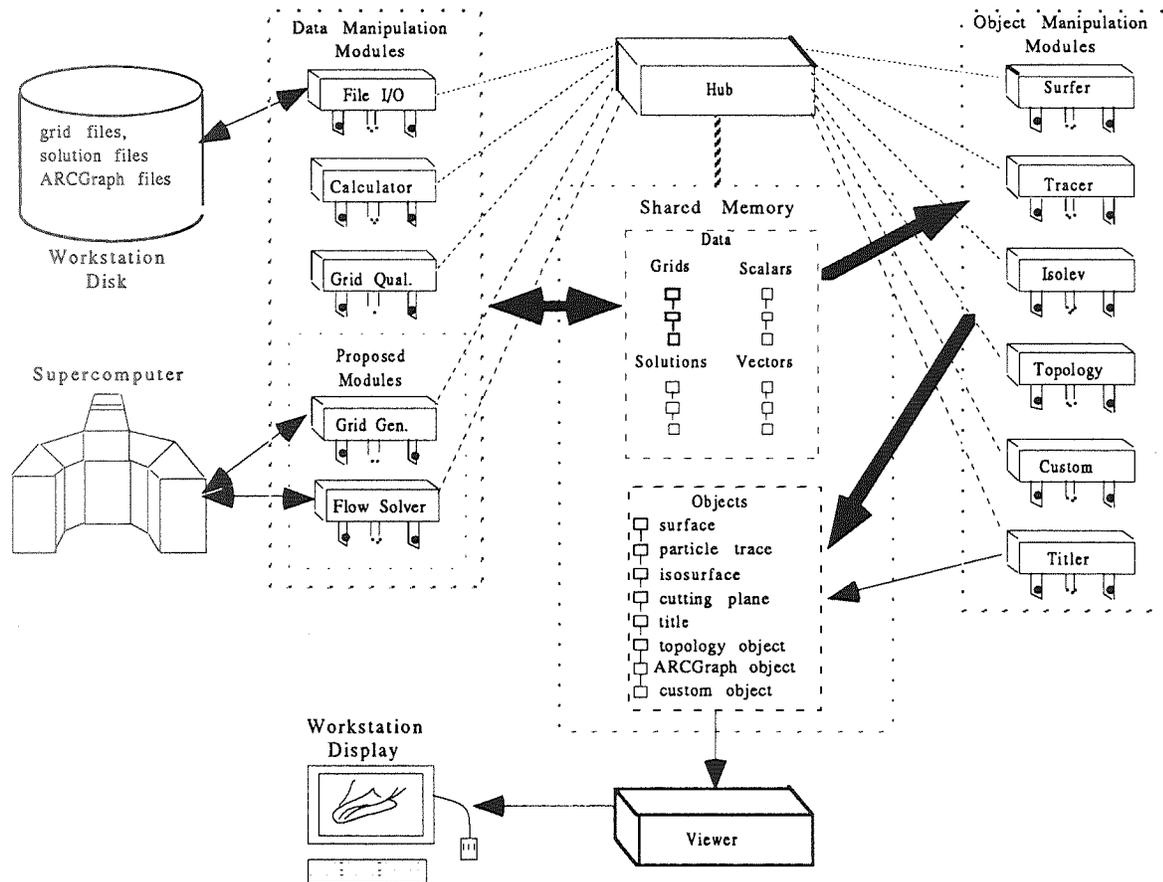


Figure 1 - FAST Architecture

CURRENT FUNCTIONALITY

In general FAST inputs grids, solutions, scalar and vector functions and creates graphical objects which are rendered to the workstation display or output to devices such as video recorders or printers. The grids describe a geometry such as an aircraft and the space around it and the solutions associate two scalar and one vector quantity with each grid node. The scalar and vector functions associate one scalar and one vector quantity with each grid node. A collection of modules operate on the data to modify it or construct graphical objects from it. The non-graphical modules make modifications to the input data which can either be written out to disk or remain in shared memory for other modules to use. The graphical modules construct graphical objects which can be interactively viewed or recorded externally to video or film media. All interaction with modules is done through a graphical user interface and each action is recorded to a script file. A FAST session can be programmed externally by constructing and executing such a script file.

Data file input and output is done with the File_io module. File_io inputs grids, solutions and function files. FAST supports all of the finite difference file formats defined by PLOT3D, an unstructured grid format developed at NASA Langley, and ARCGraph graphics metafiles. Files can be subsetted on input or output which allows for the handling of data sets too large to fit into core memory. The number of grids and solutions is limited only by the amount of core memory and swap space on the computer. All data is loaded directly into shared memory. Scalar and vector functions are loaded into registers that can be manipulated by the Calculator module.

The Calculator module operates on grids and solutions to produce scalar and vector fields that are placed in registers. The number of scalar and vector fields is limited by the number of registers (20 each) and by the amount of available

memory. There is a set of built-in CFD scalar and vector functions which are identical to the those found in PLOT3D. In addition, there is a programmable interface to the Calculator which allows users to specify formulas which apply unary and binary operators to the registers. Formulas defined by the user can be saved to and read from files.

FAST creates a variety of different surface types and rendering styles for those surfaces. The Surfer module is used to extract surfaces in computational space, which means they may or may not be planar, depending on the shape of the grid. These surfaces are planar for rectilinear data and curved for curvilinear data. Isolev, by contrast, generates either planar surfaces in physical space (cutting planes) or surfaces of constant scalar value (isosurfaces). The cutting planes can be placed at arbitrary angles to the coordinate system. Unstructured grid surfaces, cutting planes and isosurfaces can be rendered using the Surferu and Isolevu modules. All surfaces can be rendered as points, wireframe, shaded or, when displaying a scalar field, colored by a function color map. The shading and function maps can be interactively edited from the Viewer module. Vector fields can be displayed on these surfaces with or without arrowed tips. Surfaces are chosen interactively or automatically swept through the volumetric data and are illuminated by one to eight user-definable light sources.

The Tracer module in FAST is used to compute particle traces through steady state vector fields. Release points for these traces can be defined interactively and specified in either computational or physical space. Integration is either downstream, upstream or in both directions from this release point. A collection of particles (a rake) can be defined and released simultaneously. The animation of traces is used to show the dynamics of the particle paths by displaying every other integration segment (cycling) or by connecting consecutive integration segments (growing). Arbitrary scalar fields can be displayed on the particle traces by using the function color map. The Topology module also creates particle traces, but first seeks to identify critical points in the vector field by determining points where the magnitude of vector field is zero. It then classifies each critical point by analyzing the nature of the vector field in the proximity of the point, thus allowing for quick identification of interesting structures such as vortices and attachment and detachment points. Currently the particle tracing and topology features exist only for structured grids.

Traditional 2-D line plotting and text labeling is done using the Plotter and Titler modules respectively. The Plotter module can plot directly from the grid and scalar data resident in shared memory or plot coordinate pairs from an ASCII file. A variety of line styles and symbols can be applied to each plot. Multiple plots can be made on the same set of axes as to make direct comparisons of data. The Titler module provides many features to annotate a scene with scalable math and alphanumeric fonts. Text can be interactively placed in the scene or automatically aligned horizontally or vertically. Special effects like drop shadows are also available.

The Viewer module provides for the interactive viewing of the objects created by the object manipulation modules. The Viewer maintains multiple view windows and a database of scenes that can be loaded from disk or defined interactively during a session. The user is free to switch back and forth between these windows and scenes. Each view window supports full 3D interactive viewing using widget interaction, mouse or Spaceball™ input. Scenes can be antialiased and viewed in stereo using a stereo monitor. Animations can be created from these scenes and recorded frame accurately to video tape or film. FAST currently supports single frame recording to any VLAN compatible device, Abekas digital disk recorder, Lyon Lamb Minivas and Focus camera. A customizable utility is provided for recording to other devices.

Unsteady data can be visualized in FAST using the scripting language. The script file essentially runs FAST in batch mode. Objects for each time step of the visualization are created from the script, then captured on a recording device. The process is continued for each time step. A UNIX shell command feature allows for file conversions and/or network transfers after each frame is recorded.

KNOWN APPLICATIONS

The possible applications of FAST are many. The following are brief descriptions of a few of the known applications of FAST. Note that in some cases no modifications to the data or software were necessary while others required file format conversions and/or extensions to FAST.

Aerospace

FAST is used by researchers at NASA Ames and other aerospace agencies to analyze computed fluid flow around aircraft to determine the flight properties of the vehicles. Researchers use FAST as a part of their daily analysis of CFD results as well as for high-quality production of still pictures and videos for their publications and presentations.

Since FAST was originally developed for use by the aerospace community, its available visualization techniques are directly applicable. Scalar-mapped surfaces are typically used to show properties such as pressure or temperature on vehicle surfaces as well as in the free stream. Velocity vectors and particle traces indicate air flow patterns such as vortices and re-circulation regions. Other areas of interesting flow characteristics can be displayed using isosurfaces and topology extraction. Isosurfaces are used to isolate boundaries of flow characteristics. Topology extraction is used to isolate critical points in the flow field such as vortex cores and flow separation points.

Since FAST is tailored to aerospace applications there are typically little or no modifications to files or software required. Some examples of aerospace applications are shown in figures 2 and 3.

Automotive

Ford Motor Company uses FAST in their design process to investigate various characteristics around and within automobiles. Proposed designs are run through analysis codes which compute properties such as air flow and temperature distribution. The designs are potentially modified based on the results of the analysis. FAST is used to visualize the results of the analysis codes to simplify the design evaluation process.

A recent gas tank study demonstrates FAST's ability to handle both structured and unstructured data simultaneously. First, the exterior body and engine compartment were modeled using structured grids. The computed structured solution sets were then used as inlet conditions to a highly accurate unstructured model of the underbody region with emphasis on the gas tank. Using FAST, the Ford researchers were able to visually inspect and analyze their data at the points where the structured domain lined up with the unstructured regions. They were able to show the correlation between the computed results and the test data, which compared quite favorably. The visual interface that FAST offered the scientists (both in the pre-processing and post-processing phases) was most effective and essential to the analysis.

Ford's applications are possible after a few modifications. Since Ford's unstructured data format is different than that handled by FAST, a module was developed that reads the unstructured data from disk and allows the user to select the component(s) to be displayed. Ford's fluid flow data is structured in nature but not in PLOT3D format. Modifications were made to the custom module to process this alternative format. Examples of Ford's applications are shown in figure 4.

Environmental

Researchers at the National Data Processing Division (NDPD) of the U.S. Environmental Protection Agency are using FAST in several ways to visualize their computational and experimental data. Acid rain and ozone depletion models are compiled at the NDPD representing days, weeks and even years of atmospheric properties over certain regions of the U.S. These models, which are based on past history but will eventually be used to help predict atmospheric behavior, indicate definite patterns in air pollution and acid rain distribution related to weather, topology and other factors.

To visualize the acid rain models FAST is used by animating a time series of a single slice of the model representing a constant altitude level. The property being investigated is plotted on the plane as color with high values, or concentrations, of the property in bright colors and low values in darker colors. The values are sometimes accentuated by raising the points off the plane in proportion to the magnitude of the property value, thus giving the plane a relief map look. This technique makes high values easily recognizable since they are displayed as brightly colored peaks in the plot. The plots are typically animated over time with a geographical reference (eg. state outlines) and sometimes with corresponding velocity vectors representing wind direction and magnitude. The researchers will often animate the plots of one property next to another to show the correlation of the two. For

example, a time series of the occurrence of acid rain, can be displayed next to the corresponding displays of precipitation and air pollution concentration. This clearly shows that there is a correlation between the presence of both high concentrations of air pollution and precipitation with the occurrence of acid rain.

Another application of FAST at the EPA is the visualization of the changing ozone hole over the south pole. In this time series the ozone concentration levels are plotted as described above using color and height to represent the different levels. The result is a dramatic series of images showing the development of a large region of low values, or a hole, in ozone concentration. As the series, or time, continues, the hole disappears and redevelops periodically.

These visualizations were possible after converting the data files to PLOT3D format. Time-dependent data was animated by looping through sequences of 2-D planes using the Surfer module. Examples of EPA applications are shown in figure 5.

Earth/Space Sciences

At NASA Goddard Space Flight Center FAST is used by earth and space scientists to analyze their research data. In particular it was used to visualize the results of a study done on the topology of the earth's magnetosphere. In this case the researcher was interested in the effects of the interplanetary magnetic field, solar wind and ionosphere on the magnetosphere.

The interesting aspect of this visualization was the use of the particle tracing capability in FAST to identify the magnetic field lines around the earth. The general nature of the tracing algorithm allowed for the traces to be computed through the magnetic field rather than a velocity field as is typically done in aerospace applications. The path lines clearly showed the topology of the magnetic field and the general geometry of the magnetosphere. A shaded sphere representing the earth was accurately positioned in the scene to provide a frame of reference for the data.

This visualization required no modification to FAST. The data was already in the proper format since the researcher had used PLOT3D in the past. An image displaying the magnetosphere data is shown in figure 6.

Metallurgy

Research is done at Metalworking Technology to determine injection mold designs that promote even cooling of the injected material. The data from the analysis codes are a time series representing the cooling of the material in the mold. The visualization of the data with FAST allowed the researchers to easily identify areas of faster or slower cooling.

The data was in the form of a rectangular volume with geometry and temperature information at each node (as opposed to typical CFD data in which the geometry is defined by the grid). Because of this there was no way to directly identify and render the inner surface of the mold. The data was slightly modified to facilitate the generation of an isosurface that represented the inner mold. Temperature-mapped cutting planes were used to show the cooling on the inside of the mold. The temperature data on the plane was animated over time to effectively show the general cooling patterns in the mold.

As stated above, this visualization was done with minimal data file format conversion but no software modification. A single frame of the time-series is shown in figure 7.

FUTURE DIRECTIONS

Currently there are several enhancements planned for the FAST environment. The user interface will be modified to be more programmable and utilize a standard mechanism such as Motif. This change should make FAST easier to use and more portable to other workstation platforms. Modifications will be made to allow for the distribution of modules to alternate (remote) computers. This capability will allow the user to get the most out of his or her available computing environment. It should also facilitate the graphical 'tracking' and 'steering' of computational research codes running on remote supercomputers. Modifications will also be made so that the addition of modules to the FAST environment will be easier. A programmer's guide will be produced to further simplify this process. Other enhancements will include increased unstructured data visualization capability and improved unsteady, or time-dependent, data visualization.

CONCLUSIONS

It has been shown that FAST can be effective in visualizing data from several different research disciplines. With FAST being used at over 400 sites across the country there are most likely many more applications that the authors are not aware of. In the future, FAST should continue to mature as a flexible yet powerful visualization tool so that scientists in a wide variety of research areas can use it to enhance their productivity.

REFERENCES

1. G. Bancroft, et al., "FAST: A Multi-Processed Environment for Visualization of Computational Fluid Dynamics", *Technology 2001 Conference*, December 1992.
2. G. Bancroft, "Scientific Visualization in Computational Aerodynamics at NASA Ames Research Center", *IEEE Computer*, August 1989
3. P. Buning, et al., "Flow Visualization of CFD Using Graphics Workstations", AIAA 87-1180, *Proc. 8th Computational Fluid Dynamics Conf.*, June 1987.
4. A. Globus, et al., "A Tool for Visualizing the Topology of Three-Dimensional Vector Fields", *Visualization '91 Conference*, October 1992.
5. J. Helman, L. Hesselink, "Representation and Display of Vector Field Topology in Fluid Flow Data Sets", *IEEE Computer*, August 1989
6. G. Kerlick, "ISOLEV: a level surface cutting plane program for fluid flow data", *SPIE Vol. 1259 Extracting Meaning from Complex Data: Processing, Display, Interaction*, 1990
7. P. Walatka, P. Buning, "PLOT3D Users Manual Version 3.6", NASA TM101067
8. P. Walatka, et al., "FAST User Guide Version B2.1", NAS Document No. RND-92-013

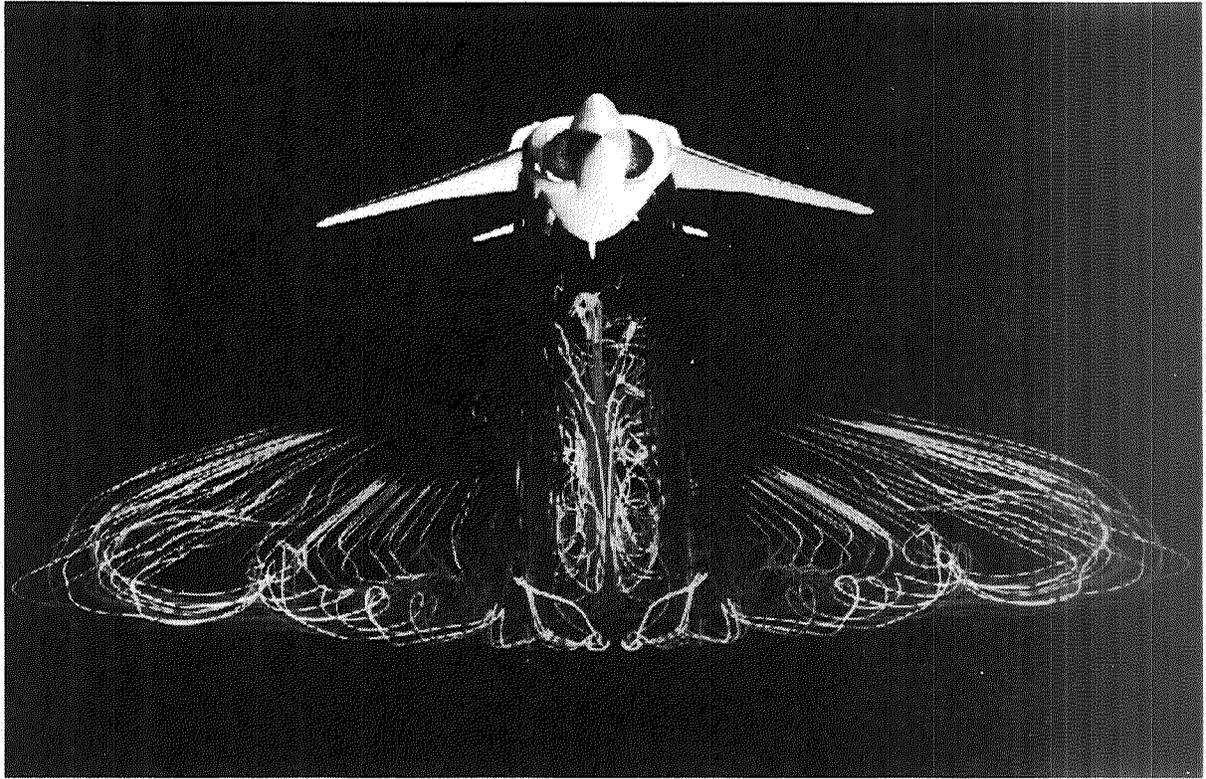


Figure 2: Harrier Jet - Data from NASA Ames Research Center

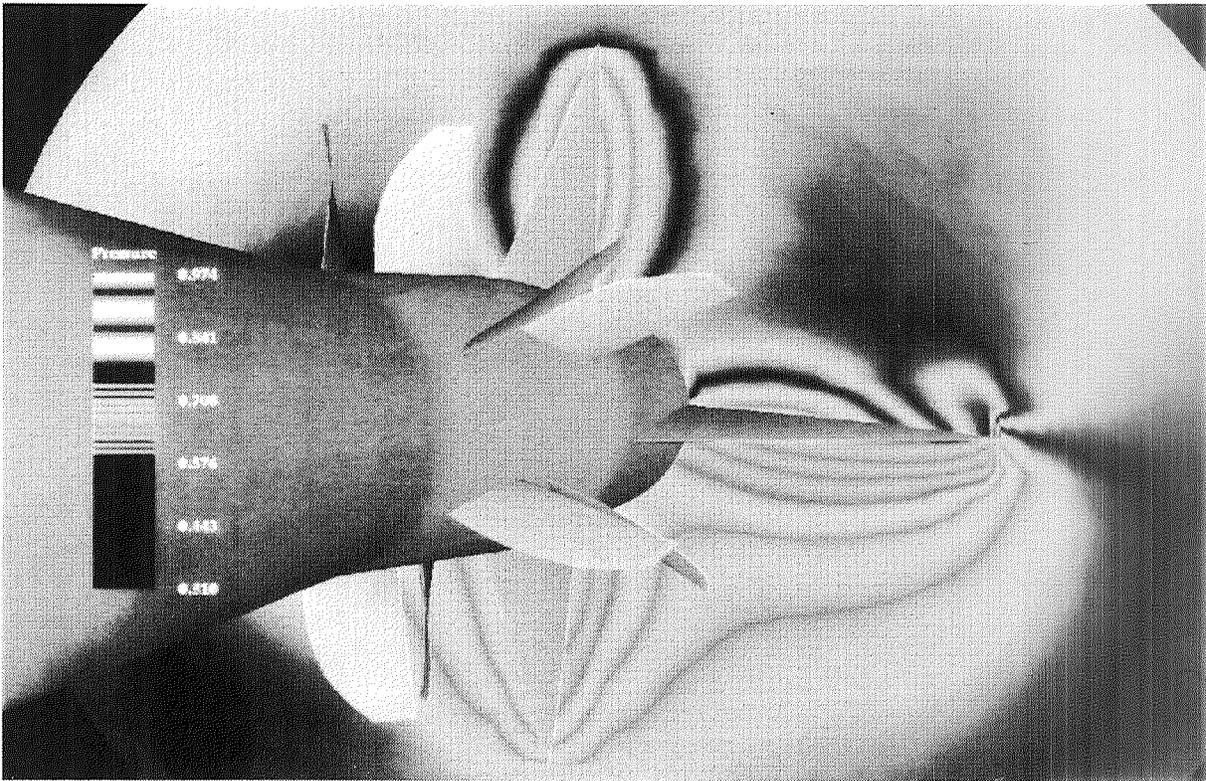


Figure 3: Cruise Missile - Data from Mississippi State University

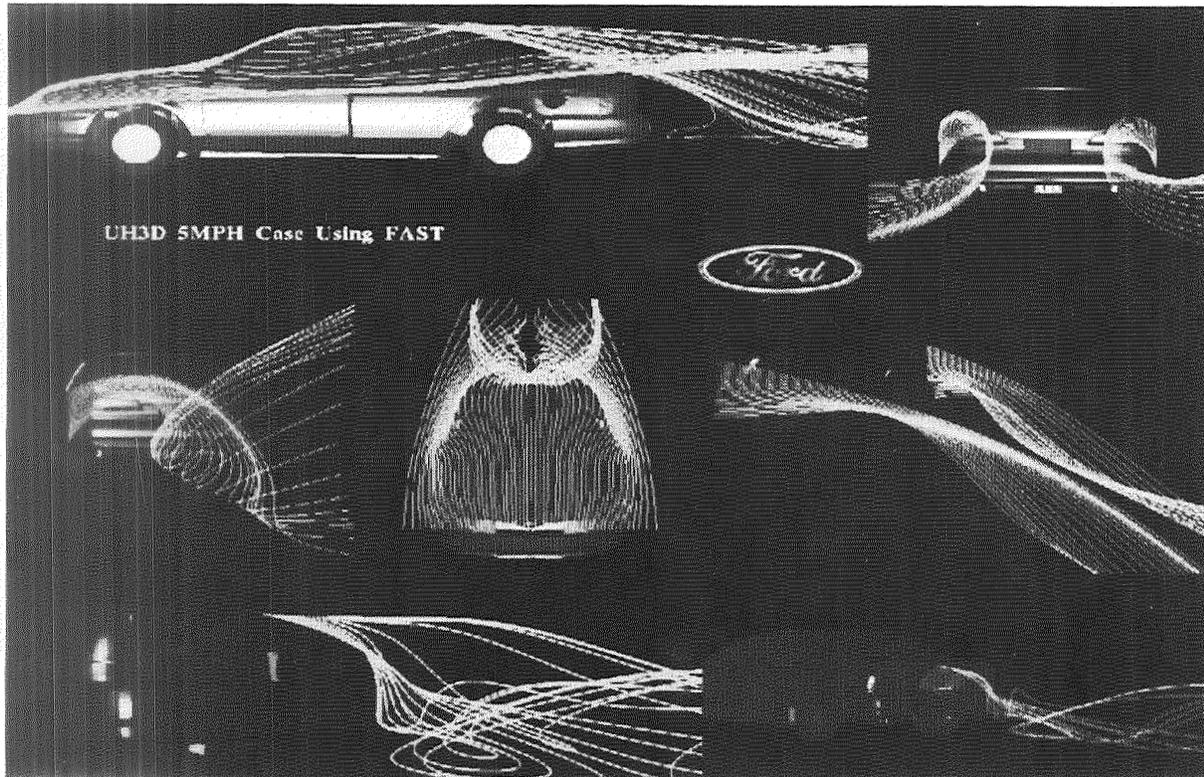


Figure 4: Taurus - Data from Ford Motor Company

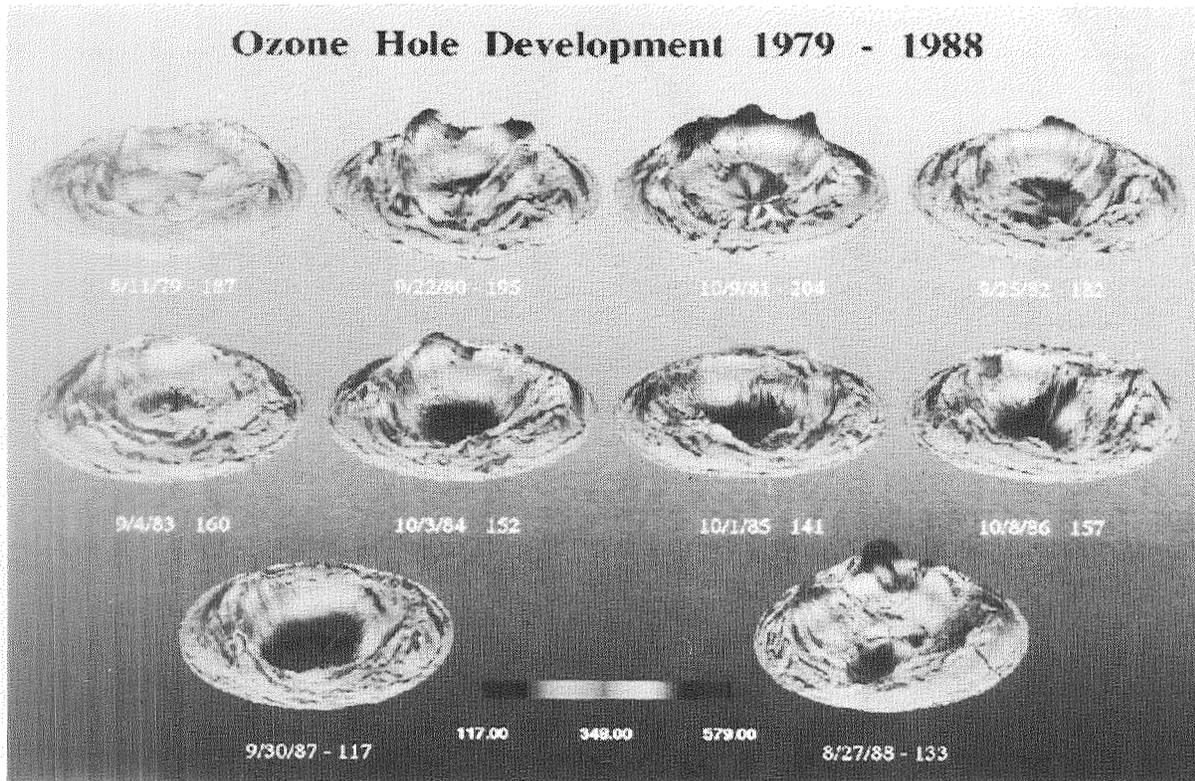


Figure 5: Ozone Hole - Data from the U.S. Environmental Protection Agency

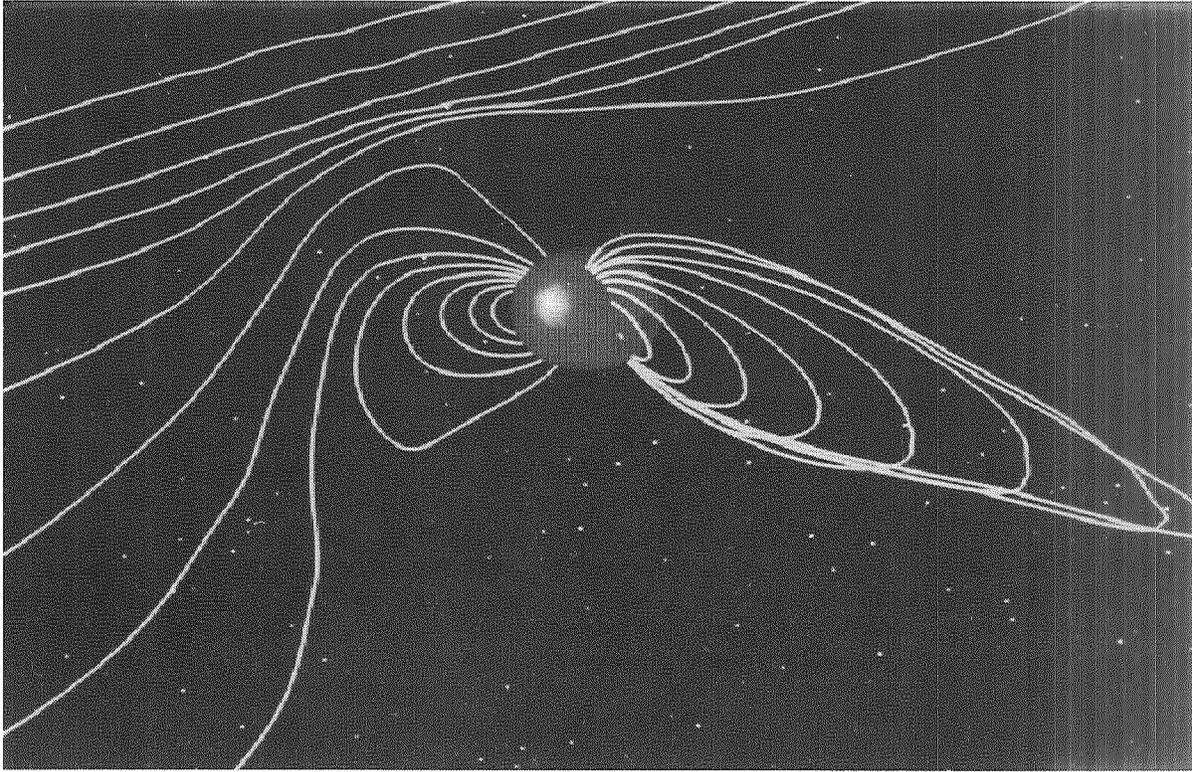


Figure 6: Earth Magnetosphere - Data from NASA Goddard Space Flight Center

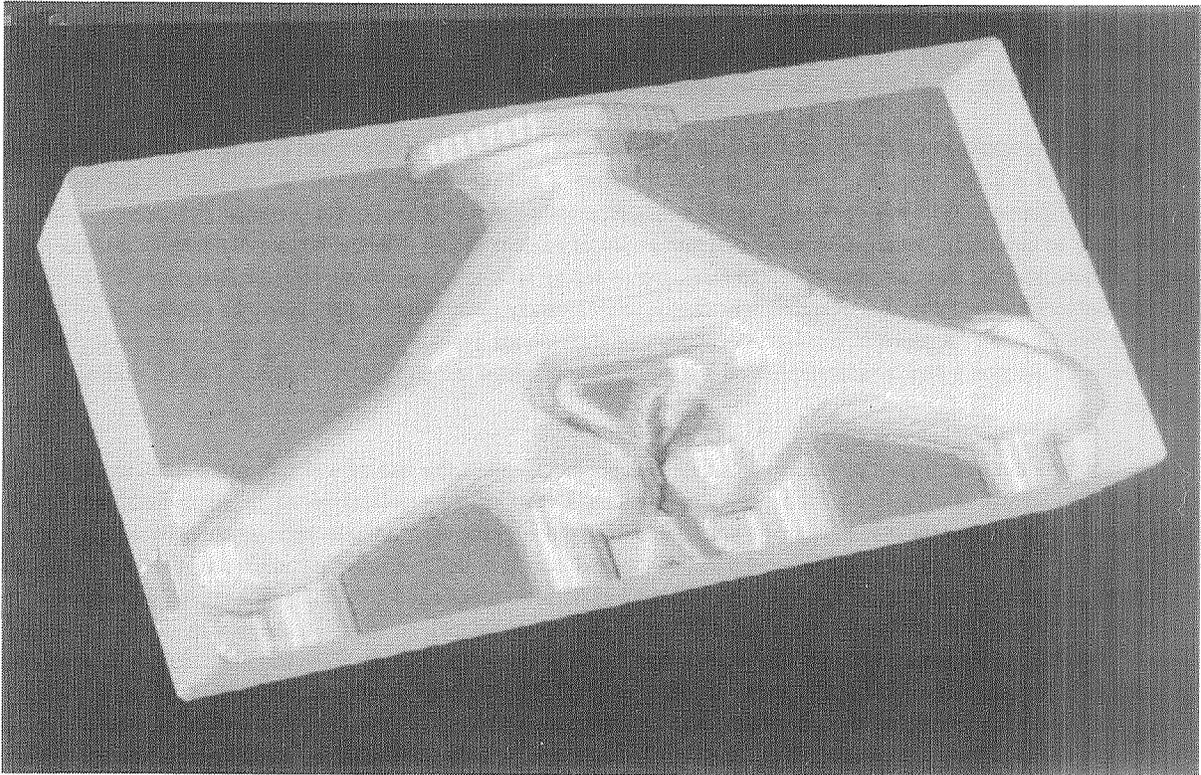


Figure 7: Injection Mold - Data from Metalworking Technology Incorporated

340-61
150510

N 93 - 25601

P-9

INTEGRATION OF DESIGN, THERMAL, STRUCTURAL AND OPTICAL ANALYSIS, INCLUDING THERMAL ANIMATION

Ruth M. Amundsen
NASA Langley Research Center
Hampton, VA 23681

ABSTRACT

In many industries there has recently been a concerted movement toward "quality management" and the issue of how to accomplish work more efficiently. Part of this effort is focused on concurrent engineering; the idea of integrating the design and analysis processes so that they are not separate, sequential processes (often involving design rework due to analytical findings) but instead form an integrated system with smooth transfers of information. Presented herein are several specific examples of concurrent engineering methods being carried out at Langley Research Center (LaRC): integration of thermal, structural and optical analyses to predict changes in optical performance based on thermal and structural effects; integration of the CAD design process with thermal and structural analyses; and integration of analysis and presentation by animating the thermal response of a system as an active color map -- a highly effective visual indication of heat flow.

INTRODUCTION

Efficiency and accuracy in analytical modeling can be increased substantially by the integration of design and analysis -- a process known as "concurrent engineering." This approach is being used in many areas to minimize the time spent in analysis without sacrificing quality. In particular, analysis effort can be minimized by utilizing the model already developed by a designer on a CAD system. Also, models and results can be transferred among analysts electronically, to avoid repetitive development of models or manual input of results.

Electronic integration of design and analysis processes was achieved and refined during the development of an optical bench for a laser-based aerospace experiment. One of the driving requirements for any complex optical system is its alignment stability under all conditions. Accurate predictions of optical bench or test bed deflections are necessary to calculate beam paths and determine optical performance. This is especially true in design of aerospace instruments, but can also apply to manufacturing or laboratory processes which use complex or high-precision optical trains. Another requirement that is increasingly demanded of any analysis process is to do it faster and better; create a more streamlined process and include all known variables to produce the best possible predictions. These goals can be accomplished by using an integrated process to accomplish design and all analyses.

The heart of the concurrent engineering process described here is the use of a single integrated model for thermal and structural analysis of a high-precision system, in this case an optical bench. The deflections of the optical bench due to structural loads can be calculated by applying the appropriate structural software package to a solid geometry model. The thermal response of the optical bench is determined by translating the same model to the thermal analysis software package. The predicted temperatures are used to calculate thermally-driven distortions. This method allows an exact calculation of the optical bench or test bed deflection due to complex thermal distributions in combination with various structural loading conditions. These calculated deflections are then used to automatically modify an optical analysis model. The change in optical performance due to given thermal and structural loads can then be determined. Designs can be optimized for peak optical performance. The analytical model can be taken directly from the design software, which eliminates an additional point for extra labor and potential error. The integrated nature of the process streamlines the modeling and analysis procedure as well as ensuring model continuity between design, thermal and structural analyses. Described herein are methods

to build an analytical model directly from the CAD model of the designed part, use the same model for structural and thermal analyses, and use the results of these analyses for the optical analysis to predict final performance.

This integrated analysis process has been built around software that was already in use by designers and analysts at NASA Langley Research Center. The PATRAN[®] solid modeling / finite element package is central to this process, since it was already in common use at LaRC. Most of the integration and interface steps described here are also possible with other packages, although certain of the translators were developed or modified for use with these specific software packages.

INTERFACE BETWEEN DESIGN AND ANALYSIS SOFTWARE

The design software currently being used in this process is Pro-Engineer[®]. A part is completely designed in Pro-Engineer, which produces a three-dimensional model of the part as well as all of the fabrication drawings. A Pro-Engineer solid shaded model of a complex assembly is shown in Figure 1. This example assembly is a laser reference cavity which is mounted on an optical bench. There are two basic methods available to translate from Pro-Engineer CAD software to the PATRAN 2.5 solid modeling software. Both of these methods have been used to produce viable models. One is to mesh the solid geometry of the part in Pro-Engineer and translate that mesh to PATRAN. The disadvantages to this method are: only the mesh is transferred, not the underlying solid geometry, so the geometry and mesh cannot be changed in PATRAN; and the mesh is limited to only tetrahedral or triangular elements. The second method is to transfer the part from Pro-Engineer to an IGES file, which is a standard graphics format, and read this file into PATRAN using the CADPAT translator. This translates the phase I (underlying) geometry, but only in the form of surfaces and lines, not PATRAN's solid geometry elements called hyperpatches. Thus the analyst must still define hyperpatches based on the geometry defined by the translated surfaces. This can actually be helpful as the analyst can choose to ignore details such as bolt holes in constructing the analytical model. The disadvantages to this method are this re-creation of the solid form from the transferred surfaces (which only applies when the part being transferred is a solid rather than plate elements), and also that during translation of an assembly of parts, the orientation of the individual parts is lost and the assembly must be reconstructed from the components.

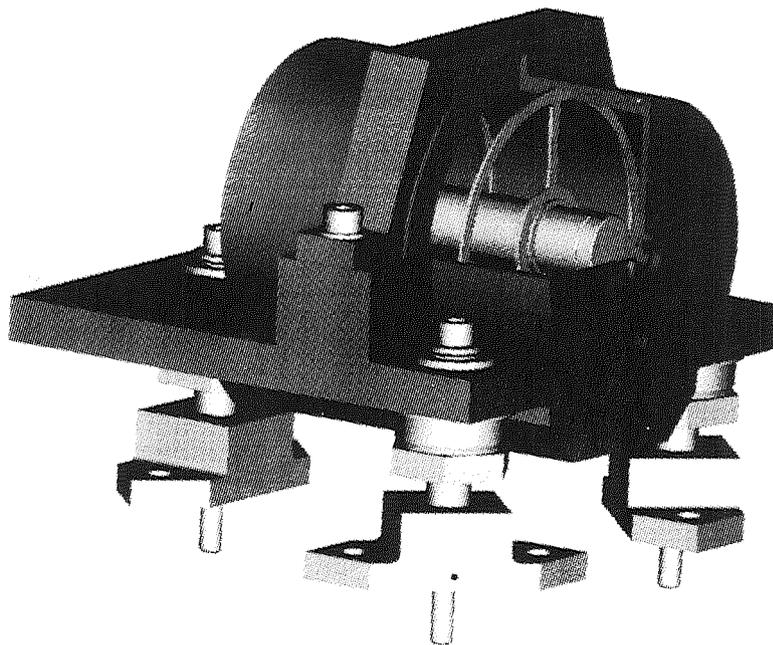


Figure 1. Laser reference cavity as designed and portrayed in Pro-Engineer

The best method in terms of the current integration process, for most parts, is to transfer a part into PATRAN by saving the 3-D part design as an IGES file. This IGES file is read into PATRAN as lines and surfaces which define the edges of the part; three-dimensional hyperpatches are then created to fill in the part. This method allows the analyst to mesh the part with the optimum element type, perform runs with different meshings, and easily alter the underlying geometry and analyze the effect of the change. Shown is an example of the mount portion of the reference cavity which was brought in as an IGES file (Figure 2) and then filled in with hyperpatches (Figure 3).

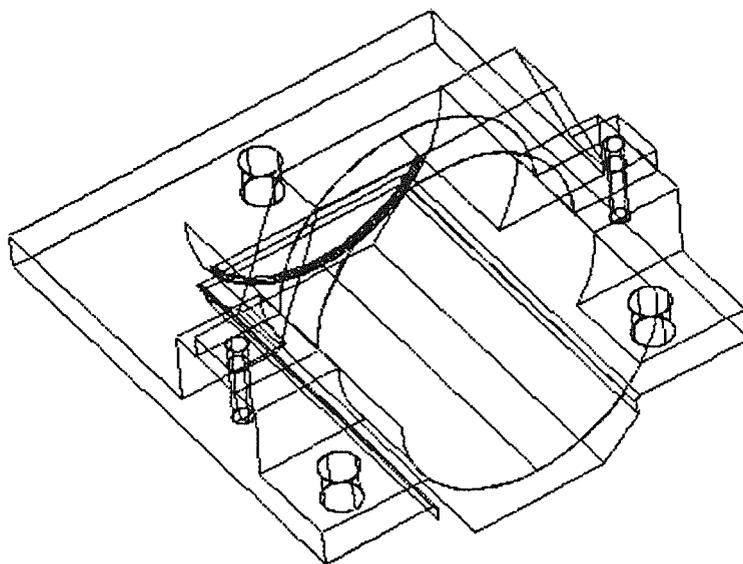
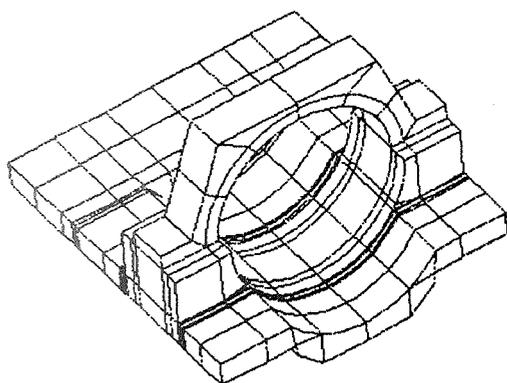
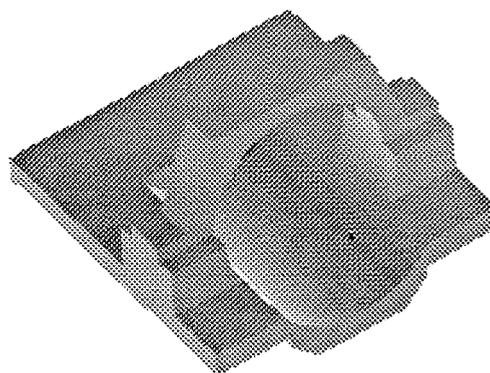


Figure 2. Mount part imported to PATRAN from ProEngineer through IGES file



(a) Shown as elements



(b) Shown as solid shaded part

Figure 3. Mount part with solid elements created in PATRAN

This design/analysis integration has several benefits. In terms of streamlining the process, there is much less work to be done by the analyst since the majority of the geometry is imported automatically. The entire process of taking dimensions from a design drawing and manually building up the geometry is eliminated. Also, the analyst is automatically working with the most current version of the design. In terms of improving the results, the fact that the human interface of re-entering the dimensions is eliminated will lessen the probability of errors in

the model. Also, geometries that are difficult to model and would perhaps be approximated are automatically translated exactly. In some cases, however, the CAD model actually has too much detail for the analyst; in these instances the model can be simplified after the transfer.

There will be an increase in capability with the releases of future versions of Pro-Engineer and PATRAN. First, a mesh of an entire assembly, rather than single parts, can be created in Pro-Engineer and sent to PATRAN. Also, the IGES file method can be used on an assembly, and the positional information of one component relative to others will be maintained. In PATRAN-3 the phase I solid geometry will not be required, so that the analyst could create elements directly from the imported CAD geometry. PDA Engineering will be adding an option to allow the underlying solid geometry to be imported from Pro-Engineer to one of the later releases of PATRAN-3. This would mean that the entire solid model, including hyperpatches, would be available to the analyst in PATRAN. As the interface improves with future versions of the programs, it will also be possible to automatically import and export changes to the design that are made either by the designer or by the analyst. These benefits would be even more striking in an industrial process, where many more designs are produced and analyzed than is common in the aerospace industry.

INTERFACE BETWEEN THERMAL AND STRUCTURAL ANALYSIS

The translations between structural and thermal analytical models are already built into the PATRAN system. However, there are a few methods that make this type of translation easier and more effective. The model can be built in PATRAN by either analyst; however there must be communication between the analysts before the model is built, so that the final model will have a structure and level of detail appropriate for both analyses. One unique aspect of the work described herein is that the structural and thermal analysts determined together what method would be best for both of them in modeling certain parts, before the model was developed. A requirement on the thermal side, which must be maintained in the model in order for it to be useful for the thermal analyst, is that between every pair of connected elements, all corner nodes must be identical. Also, the best translation to a thermal model is currently achieved by using solid elements rather than plate elements in most cases. Many of the connections between solids and plates, and plates-to-plates, that are correct for the structural analyst, do not work correctly for the thermal analyst. In order for each analyst to be able to easily create their own mesh, or use the same mesh, the phase I geometry must meet the requirements of both analysts.

The PATRAN model is translated to SINDA-85, a finite difference thermal analyzer, using the PATSIN translator. This SINDA-85 model is used to perform thermal analysis, with some modifications such as adding power sources. The structural analysis can be performed either in P/FEA (a software package that directly interfaces with PATRAN) or after translation to NASTRAN or EAL (Engineering Analysis Language). The analysts sometimes desire different levels of detail; thermal analysis commonly uses a lower level of detail than structural analysis. In that case, an identical PATRAN phase I geometry of patches and hyperpatches is still used and each analyst can create their own mesh. The thermal results from SINDA can be translated using the SINPAT translator into temperature data files which are read by PATRAN. These temperatures can be used to impose accurate thermal loads on the structural model regardless of whether the meshing is the same, as long as the phase I geometry has not been changed. Shown in Figure 4 are examples of models which were built and meshed by one analyst and used by both. This has also been done with a model using two different meshings and element numbering schemes; the interpolated values imported from one model to another were checked and found to be correct.

The only change that must be made to alter the model between use by the thermal and structural analysts is a re-definition of the material properties, usually a five to ten minute task. The material identification is maintained through the transfer; only the actual material properties need be re-input. Unfortunately the material properties are exclusive, so that each time the PATRAN model is transferred between analysts, the material properties must be redefined. This is not normally a problem since the transfer is done only once. Also, improvements slated for PATRAN version 3 will do away with this concern.

The easiest way to use the nodes and conductors created by PATSIN is to separate them into files which are called into the SINDA model using an INCLUDE statement. Thus the SINDA model can contain other data such

as heating arrays; if there is a change to the PATRAN model it will only affect the included files, with the main SINDA model left unchanged. Also, the output of PATSIN is often quite bulky, which would make editing of the full SINDA model more difficult. Using included files limits the size of the SINDA model file, and allows several different SINDA files to reference the same node and conductor files.

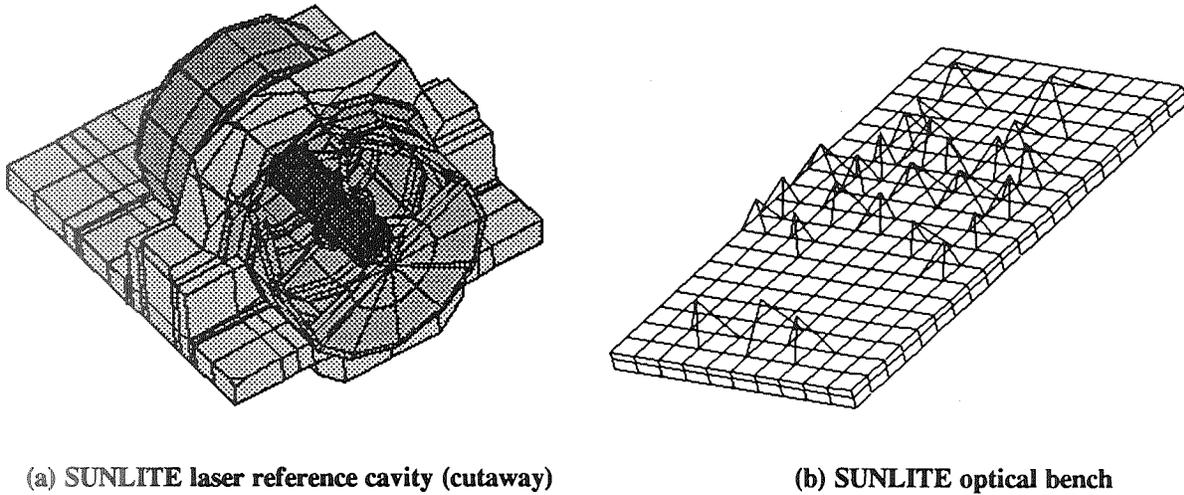


Figure 4. PATRAN models used for both thermal and structural analysis

The thermal results, either from a steady state analysis or from time steps in a transient run, are saved in a text output file. This file is operated on by the translator SINPAT to produce element and nodal temperature files that can be read by PATRAN. These files can be read directly into PATRAN, and the thermal results mapped onto the model geometry. This is shown in Figure 5 with a thermal map displayed on the laser reference cavity.

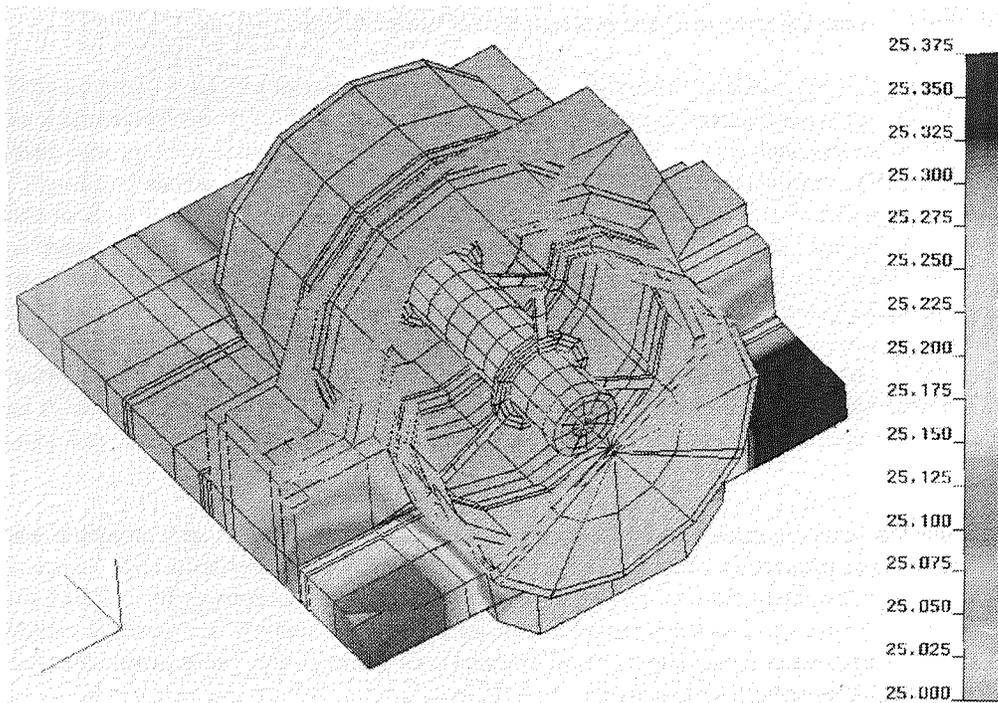


Figure 5. Map of thermal distribution (in degrees C) on laser reference cavity in PATRAN

In order to use the nodal temperatures as actual thermal loads rather than only for display, the files must be run through a program called **READER** which translates the text files to binary. The results can be interpolated onto the structural model using a built-in utility of **PATRAN** (**TEMP**, **ADD/INT**). The thermal results, imported into the **PATRAN** model, can be used in the structural analysis software to calculate thermally-driven stresses and deflections based on the predicted temperature distribution. These thermal stresses can be summed with any load-driven stresses, to produce a total reaction of the system to the environmental constraints.

INTERFACE BETWEEN STRUCTURAL AND OPTICAL ANALYSIS

Most optical models start with the assumption that the system is aligned and at rest. The optical analyst inputs surfaces, sources and objects at their designed location, and determines the performance of the system. The optical code currently used by analysts at LaRC is **CODE-V**[®]. During actual operation of the optical system, there will often be factors that cause distortions to the perfectly aligned system. In the case of an optical bench which has optical components mounted on it, there can be thermal gradients across the bench which will cause minute warping of the bench, but which result in significant distortion of the optical system from its baseline aligned performance. There can also be structural loads imposed which cause deflections, and both the thermal and mechanical loading environment can be changing with time. There is an existing translator that will look at the deformation of a single optical element such as a lens in **NASTRAN**, and translate the appropriate information to **CODE-V** to determine the distorted lens performance. However, for the optical bench structure, a method was needed to look at changes in the overall performance based on distortions of the entire bench, not only a single element.

The method that has been developed starts with writing an output file of nodal deflections from the structural analysis software. The deflections can be due to thermal, structural or any other loading conditions. A relational file is developed for that model which relates the nodes in the **PATRAN** model to the optical surfaces in the **CODE-V** model. This relational file can be used for any translations of results from that **PATRAN** model to **CODE-V**. Translation software (temporarily named **BENCH**) was developed to read the deflection file, the relational file, and a copy of the undeflected **CODE-V** model. It produces a new **CODE-V** model which has new positions and angles for the optical elements based on summing the predicted deflections and the original positions of the elements. It also adds a title to the **CODE-V** model based on the deflection case that was run and augmented by user input. For a single optical bench there is only one **PATRAN** model, but there can be a separate **CODE-V** model for each optical path. The translation must be run for each optical path for which deflection analysis is desired. **CODE-V** can then be run on the new model, and optical performance based on the distorted bench is predicted. The translation can be run for a series of time steps, using deflection results files for each time step, to predict the performance of the system as a function of time. Figure 6 demonstrates the steps of this process pictorially by showing the temperature distribution on an optical bench (a), a map of the thermally-driven distortions from this distribution (b), a map of distortions due to structural loading (c), and a map of the combined distortions (d). These combined distortions, consisting of six numbers for each optical surface (rotations and translations in each of three axes), were written to the deflection results file. This file was used to modify the **CODE-V** optical model, and yielded the predicted performance of the system under these deflections. Development of a users' manual for this translation software is currently underway.

THERMAL ANIMATION

Structural analysts commonly use animation in their presentation of results. Animation of mode shapes or predicted deflection patterns is a vivid method of capturing and conveying all the necessary information. This is done less often with thermal analysis, with the result that many viewers have a less concrete idea of the physical progression of temperatures or heat flows. A visual animation of the thermal map, in color, gives a very effective representation of the physical transfer of heat.

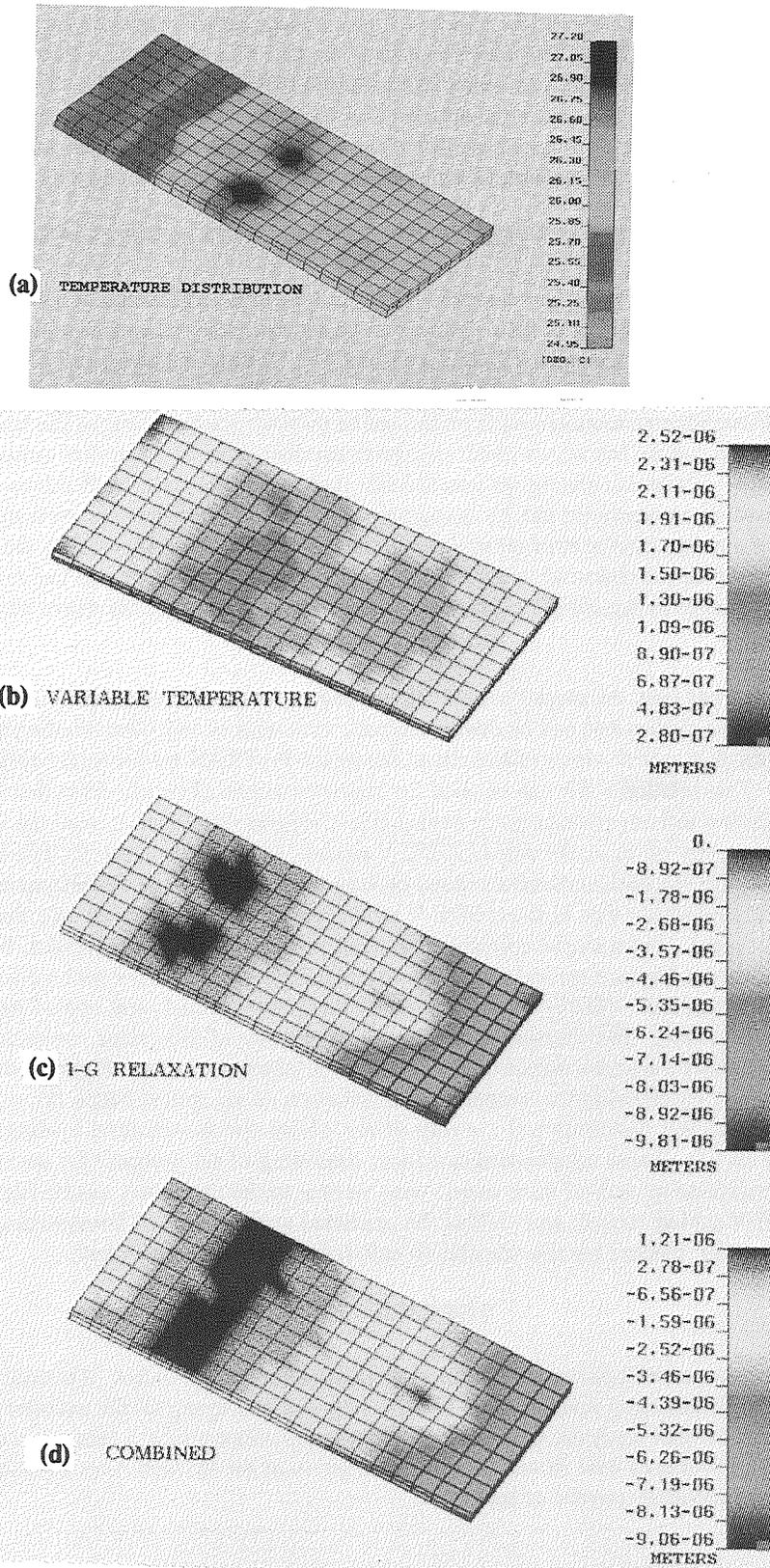


Figure 6. SUNLITE optical bench with progressive mappings: (a) thermal distribution, (b) thermally-driven deflections, (c) mechanical load-driven deflections, and (d) summed deflections

Animation of transient thermal results, in combination with an integrated structural-thermal model, is a very effective tool that has been utilized through PATRAN and its connection with SINDA. The temperatures of a part are mapped onto the geometry using a color scale. Color maps are generated for several sequential time steps. The mappings are viewed in a sequence that is run repeatedly on the screen. The progression of temperatures along the part as a function of time is observed as an animated color thermal map. This function is invaluable when evaluating the driving force behind a given reaction, and can also give an audience a much clearer understanding of the processes involved in a complex reaction. Cases can be recorded on video tape and used to demonstrate results to a larger audience. This function is also quite valuable to the analyst, as it provides a method for debugging the model and perhaps finding errors that would be time-consuming and tedious to find in any other way.

The specific method for performing this animation is to run SINDA with thermal output at all desired time steps. A virtual temperature (VTE) file must still exist from the PATSIN translation for the SINPAT translator to use in calculating nodal temperatures from element temperatures. Once SINPAT has created the nodal results files, the easiest way to set up an animation is to set up a PATRAN session file which reads in all of the frames. This avoids manual keying of three input lines per frame. The session file sets up the number of frames and the spectrum to be used. For each data frame, RUN, CONTOUR and RUN, HIDE commands are performed. The last line for the last frame must be typed in manually. After the last line is entered, the animation will begin running on the screen. The animation characteristics may be altered in real time using the animation menu. On 3-D workstations the part can be rotated on the screen during animation, so that thermal progressions on all sides of the part may be viewed.

FUTURE PLANS

Several improvements have been planned for this process in the future. The translation of bench distortions from NASTRAN to CODE-V could be integrated with the built-in capability of CODE-V to accept the distortions of a single optical element from NASTRAN to produce a complete prediction of the system performance. Also, a method could be developed to run all of the optical path translations in a model at one time. The entire process, from design software through thermal/structural analysis to optical performance, should be run on a simple system design to evaluate the efficiency of the entire technique.

The advent of PATRAN-3 should add many valuable features to the process, in terms of simplifying the CAD/analysis interface, upgrading the structural analysis flexibility, and improving the post-processing and display capabilities. Also, the greatly improved user-friendliness of version 3 may encourage more thermal analysts and CAD designers to integrate their work with the structural models. However, there may be a temporary set-back in that PATRAN-3 will probably not initially support the PATSIN translator. This lack of thermal support within version 3 will force the thermal analysts to continue using version 2.5 for at least some portion of their work, while the structural analysts shift to using version 3. It is not yet clear how much of an impact this will have on the integrated process, and how rapidly PDA will implement the SINDA interface in PATRAN-3.

CONCLUSIONS

The process that has been developed uses existing software and translators, along with some newly developed translators, to electronically link design with analysis, and to integrate several types of analysis. This allows an analyst to quickly develop a complex geometric model with much less time-consuming and repetitive labor, since the geometry is imported electronically. Since both the thermal and structural analyst can use the same model, the labor involved in creating an extra model is eliminated. With the use of an integrated model, the thermal gradients can be directly understood in terms of deflections, and presented with the other structural results. Transfer to the optical analysis model allows interpretation of the impact of worst-case environmental conditions on the optical performance. Without this, optical performance must be roughly estimated or calculated by manual input of worst-case deflections. Visual observation of the thermal dynamics of a system has been very useful in understanding and presenting the significant thermal forces in a system. This allows the analyst to concentrate

concentrate detail in the appropriate portions of the model and uncover errors that would otherwise be undetectable. All these improvements lead to better, faster and cheaper accomplishment of total system design.

Many of the described aspects of integration and animation have high potential for commercial and industrial applications. The integration between design and analysis can be used to streamline any mass-production application such as design and manufacture of automotive parts, machine equipment, or plastics fabrication. The integration between thermal, structural and optical analysis can be useful in any field which requires maintenance of close tolerances between optical components; examples are automated fabrication/assembly lines which use lasers for position measurement, scientific laboratories which use lasers for experimentation and measurement, and of course research and development centers such as LaRC in developing sensitive optical instruments for applications such as remote sensing of pollutants. The animation of thermal analyses can be of significant use in any field where understanding of thermal flow is critical. Examples include plastics manufacturing (such as molding and extrusion processes), automotive engine design, electronic design and fabrication, analysis of chemical reactions, and power plant design.

ACKNOWLEDGEMENTS

The efforts of Kelly Smith in structural analysis, Steve Hughes in design, Maria Mitchum in software development, and Greg Herman, Alan Little and Andrew Cheng in optical analysis are gratefully acknowledged. The funding for this work was provided by the SUNLITE project.

DATA SYSTEMS DYNAMIC SIMULATOR

Christopher Rouff
Melana Clark
Bill Davenport
NASA Goddard Space Flight Center
Code 522.1
Greenbelt, MD 20771

Philip Message
Stanford Telecom
7501 Forbes Boulevard
Seabrook, MD 20706

541-61
150571
N93-25602
p-9

ABSTRACT

The Data System Dynamic Simulator (DSDS) is a discrete event simulation tool. It was developed for NASA for the specific purpose of evaluating candidate architectures for data systems of the Space Station era. DSDS provides three methods for meeting this requirement. First, the user has access to a library of standard pre-programmed elements. These elements represent tailorable components of NASA data systems and can be connected in any logical manner. Secondly, DSDS supports the development of additional elements. This allows the more sophisticated DSDS user the option of extending the standard element set. Thirdly, DSDS supports the use of data streams simulation. Data streams is the name given to a technique that ignores packet boundaries, but is sensitive to rate changes. Because rate changes are rare compared to packet arrivals in a typical NASA data system, data stream simulations require a fraction of the CPU run time. Additionally, the data stream technique is considerably more accurate than another commonly-used optimization technique.

INTRODUCTION

Development of the Data Systems Dynamic Simulator (DSDS) started in the late 1970's at Marshall Space Flight Center. Under contract to NASA, the General Electric Company was tasked to build a discrete event simulation tool especially suited for modeling NASA end-to-end data systems of the Space Shuttle and Space Station eras. Since then, DSDS has been in continual use. In 1985 the management and control of DSDS was transferred to Goddard Space Flight Center.

In 1986, Stanford Telecom was tasked by the Customer Data Operations Service (CDOS) project at GSFC to develop an end-to-end model of CDOS. Although CDOS, now called EDOS, is a ground-based system, it was necessary to accurately model the sequence of data arriving at EDOS (from space). The objective of the EDOS model was to study the traffic and mission profile for selected twenty four hour periods. As in any modeling task, the proper tool selection is an essential step. Therefore, EDOS personnel searched for a suitable tool to meet the EDOS model requirements.

Because of EDOS' high data rates, complexity and excessive simulated time requirement, the search proved to be difficult. None of the tools or languages surveyed supported a suitable optimization technique. It was decided that the data stream modeling optimization technique offered the best potential solution. While the data stream theory was well known, the method had never been implemented on a scale suitable for the EDOS model. A survey of existing tools was done, and DSDS was selected for the following reasons: it was NASA-owned and maintained, it was extensible and it supported orbital modeling.

Since 1986, data streams have been used to model other NASA end-to-end data systems, including Space Station Freedom (SSF) and the Earth Observing System (EOS).

DSDS OVERVIEW

DSDS contains predefined, configurable elements that can be used to represent components of a data system. Examples of elements that are available for simulating systems are CPU's, data generators (scientific instruments or experiments), orbit calculators and schedulers. A model is constructed by linking elements together to represent a network over which simulated data will flow. Sizes of queues and time to send data between elements can be simulated. From the results of the simulation the sizes of buffers, the speed of processors and the speeds of data links that will be needed for the system to have a particular performance can be determined. Though DSDS itself does not process cost information, once the types of components are determined, the modeler can then determine the price of the proposed system.

The model assumes that data is sent between elements in a packet. The packets can be simulated in one of two modes: packet by packet or as a data stream. Using packet simulation, packets are sent through the system one at a time. In the data streams model the system is modeled by using the rate at which data is being sent through the system. Both simulation techniques are discussed further below.

As shown in Figure 1, there are four main parts to DSDS: SETUP, SIMULATE, GRAPHICS and USRLOAD. SETUP is where the user defines the elements to be used and how they are to be interconnected. Parameters for each element allow the elements' behavior to be customized for different applications. SIMULATE is where a model of a system is simulated. SIMULATE displays graphics on the progress of the simulation, displays other statistics, has debugging capabilities to trace through a simulation, and produces reports that provide statistics on queue lengths and time for data to flow through the system. GRAPHICS prints plots and reports generated by SIMULATE. USRLOAD lets users extend DSDS by allowing them to write their own elements.

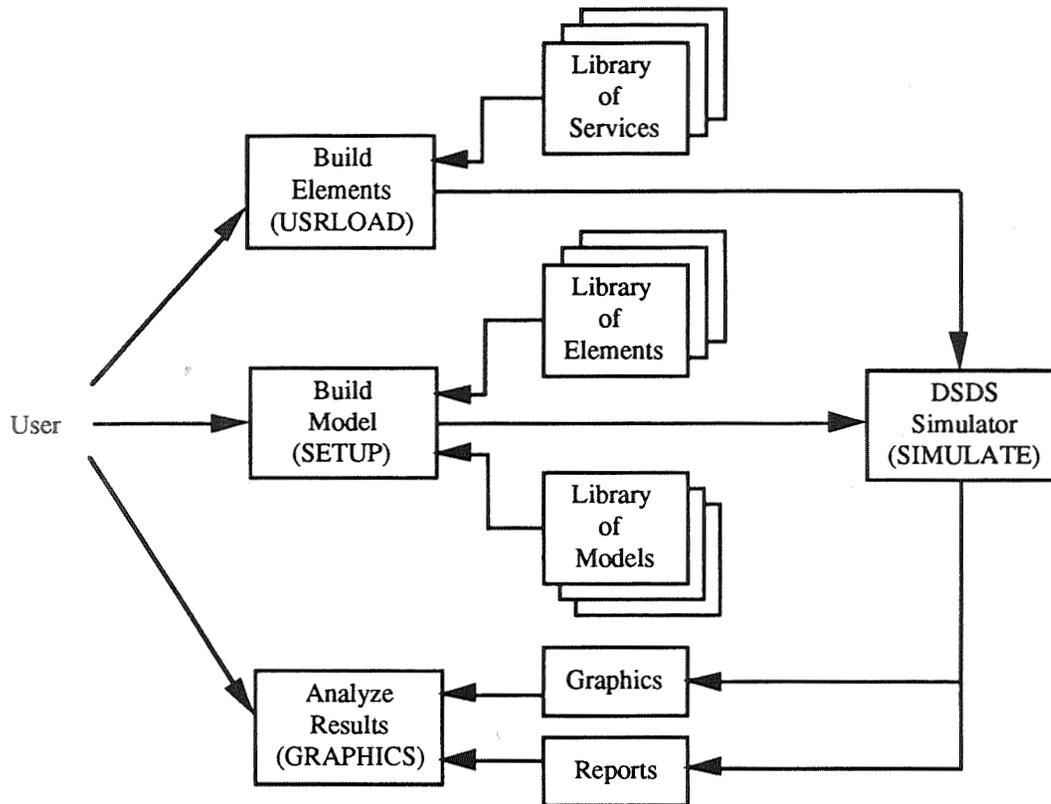


Figure 1: Components and flow of data in DSDS.

DSDS is a block-oriented simulation tool. Standard DSDS consists of a library of about fifty pre-programmed blocks (i.e., elements). The use of standard DSDS blocks require no programming. Each

DSDS block has a unique number identifier plus a name. Block names identify a type of operation such as CPU, DUPLIC, AND, OR, BRANCH, etc. Alias names are optional.

DSDS users are prompted to choose default values or enter appropriate parameters such as packet size, priority, data type, destination block, processing rate, etc. DSDS elements are connected by lines or links that are specified by the model builder. During model execution, information is passed from block to block via these links.

In DSDS, the movement of information is called an event. A DSDS event is represented by a twenty five word FORTRAN array. Array words contain specific attributes about the event (e.g., name, priority, destination, number of bits, data type). A FORTRAN subroutine is associated with each block. The arrival of an event invokes the subroutine code which in turn performs the simulated operation.

Each event arrival changes the state of the system being modeled. System states are recorded in a time-ordered history, called a timeline file, that is created and maintained by DSDS as the simulation proceeds. Data from a timeline can be converted interactively or off line into customized graphs and/or reports. Events occur in simulated time. The next event may happen in a micro-second or a day. Therefore there is no relationship between wall clock time and simulated time. However, the number of events processed has a direct relationship to CPU run time. This issue will be discussed later.

The user of standard DSDS is constrained by the functionality of a finite set of blocks (or elements). Standard DSDS constraints are indicative of similar tools with pre-programmed elements. The constraints are twofold. The first is the nuisance factor. While it may be possible to get the job done with standard elements, the model configuration is not economical. In these cases, simplification could be accomplished with a customized element. In the second, more serious case, functionality does not exist. Then it becomes essential that a new element be constructed.

USRLOAD supports the development of user-defined elements. The modeler is provided with a template and prompted to describe characteristics, such as name, number of user parameters, number of queues, etc. A documented set of DSDS utility functions is available to the user. The use of these utilities guarantees that the user-defined element will function consistently with the elements in the standard set. However, the user is required to write a FORTRAN subroutine with the desired functionality.

The integration of user-defined elements into the standard set is decided by the DSDS manager. Examples of user elements which have been ported to the standard set include two TP-4 protocol elements. These elements were developed by a DSDS user (LinCom Corporation) for Johnson Space Flight Center (JSC). The user extension feature gives rise to the term "Dynamic" in the name Data Systems Dynamic Simulator.

OPTIMIZATION TECHNIQUES

Large complex data systems create an extremely large volume of data. A limiting factor, when modeling these data systems, is the amount of simulation (CPU) run time required to complete the simulation. The CPU run time increases when the number of events in the simulation increases. When generated each packet creates an event. As the packets flow through the model (being queued, processed, sorted, etc.) more events are created, thus increasing the CPU run time of the simulation. For large complex data systems with extremely large volumes of data the number of events generated can become overwhelming. Since the purpose of doing a simulation is to try different configurations, determine potential bottlenecks, and test out the effects of different components with different costs and performance, it is important to be able to run a simulation many times in a timely fashion. This allows different versions of the model to be compared to determine which components would work best given system constraints, such as price and performance.

To bring simulation run times down to a reasonable value, modelers often use optimization techniques. The next two sections describe two different optimization techniques and their use with DSDS.

Artificially-Inflated Packets

When modeling a NASA system, the data that passes through it is divided into packets. Packets (or messages) may be characterized by size. In standard DSDS, packet size is a parameter that is supplied by

the user. The size of these data packets can vary, but the typical NASA packet size is in the range of 10,000 bits. Very large numbers of packets are associated with the simulation of current and emerging NASA data systems. Every packet transfer results in a state change within the model. Each state change requires CPU run time. For instance, to simulate a 24 hour "true" EDOS model it would take more than a day of CPU time on a relatively fast computer (VAX 8600).

As complexity, data rates and need for longer model run times increases, so does the need for an optimization technique. Artificially raising the packet size is one way of optimizing the simulation. By increasing the packet size, say from 10,000 to 100,000 bits, CPU run time is reduced by a factor of 10 due to the decrease in events flowing through the system. Unfortunately, increasing the size of packets results in errors. The magnitude of the error can be predicted by comparing the results to a "truth" model. A truth model is one constructed with the "real" packet size. It is compared to a test model which is constructed with the elevated packet size. However, it is impractical to develop a truth model for some large systems, such as EDOS. The model, when fully configured, would consist of more than 100 payloads (experiments). Before reaching users, payload data must pass through 10 or more processing points. Thus it is not easy to determine the effects on the fidelity of the modeling results when the packets are inflated artificially.

Data Streams

The data streams technique is another way of optimizing the modeling process. As stated above, the objective is to reduce the number of events in a given simulation run. The data streams method models rate changes rather than individual packets. Consider the example of an experiment which transmits a 10,000 bit packet at the constant rate of 100 packets per second. Consider also that the experiment transmits 10 minutes each orbit. In a "true" model this translates into 600,000 packets (events) per 10 minute duty cycle. A data stream represents this duty cycle as two events; namely one start and one end. (The data stream would be characterized as a 10 minute stream with a transmission rate of 600,000 bits per second.)

The key to understanding the data stream methodology is that it takes advantage of the linear flow of data between state changes. The speed at which the simulation runs for the data streams method is not dependent on the volume of packets, but instead is dependent on the number of times the data rate changes. The data streams method requires less computation and thus reduces the time to simulate the passing of data through the system. Data streams take advantage of the fact that data systems behave linearly between state changes. Therefore, data streams can model the effects of the changes in the data rates of a system rather than modeling each individual packet. This optimizing method will also reduce the CPU run time of a simulation due to the decrease in events.

One difference between packet and data stream modeling is the way a processor's bandwidth is allocated. A packet, no matter its size, occupies the entire bandwidth of its processor for some finite period of time. Data streams, on the other hand, occupy only that portion of the bandwidth which is equal to or less than its transmission rate. Furthermore, data streams share the bandwidth proportionally on a first in, first out basis with other competing streams.

Comparison of Data Streams to Packet Modeling

Data stream simulation is tantamount to modeling with infinitesimally small (approximately 1 bit) packets. In terms of magnitude, 1 bit is closer to 10,000 than 20,000 bits and more. Based upon this empirical point alone, an assumption could be made that errors induced by data stream optimization would be less.

An experiment was performed to assess the error produced by artificially increasing the packet size of a packet model of a system and the error introduced by data streams modeling. Models using these optimization techniques were compared to a "Truth Model" - a model which is run using the actual packet size. Figure 2 shows the model that was used for the experiment. Figure 3 shows the run times of the Truth Model where the packet sizes were 15 Kilobits (reflecting the actual implementation), a model where the packet size was artificially expanded to 1 Megabit packets, and the data streams methodology. As expected, the Truth Model ran the longest, 3,279 CPU seconds. The expanded packet model ran for 47 seconds and the data streams model for 62 seconds.

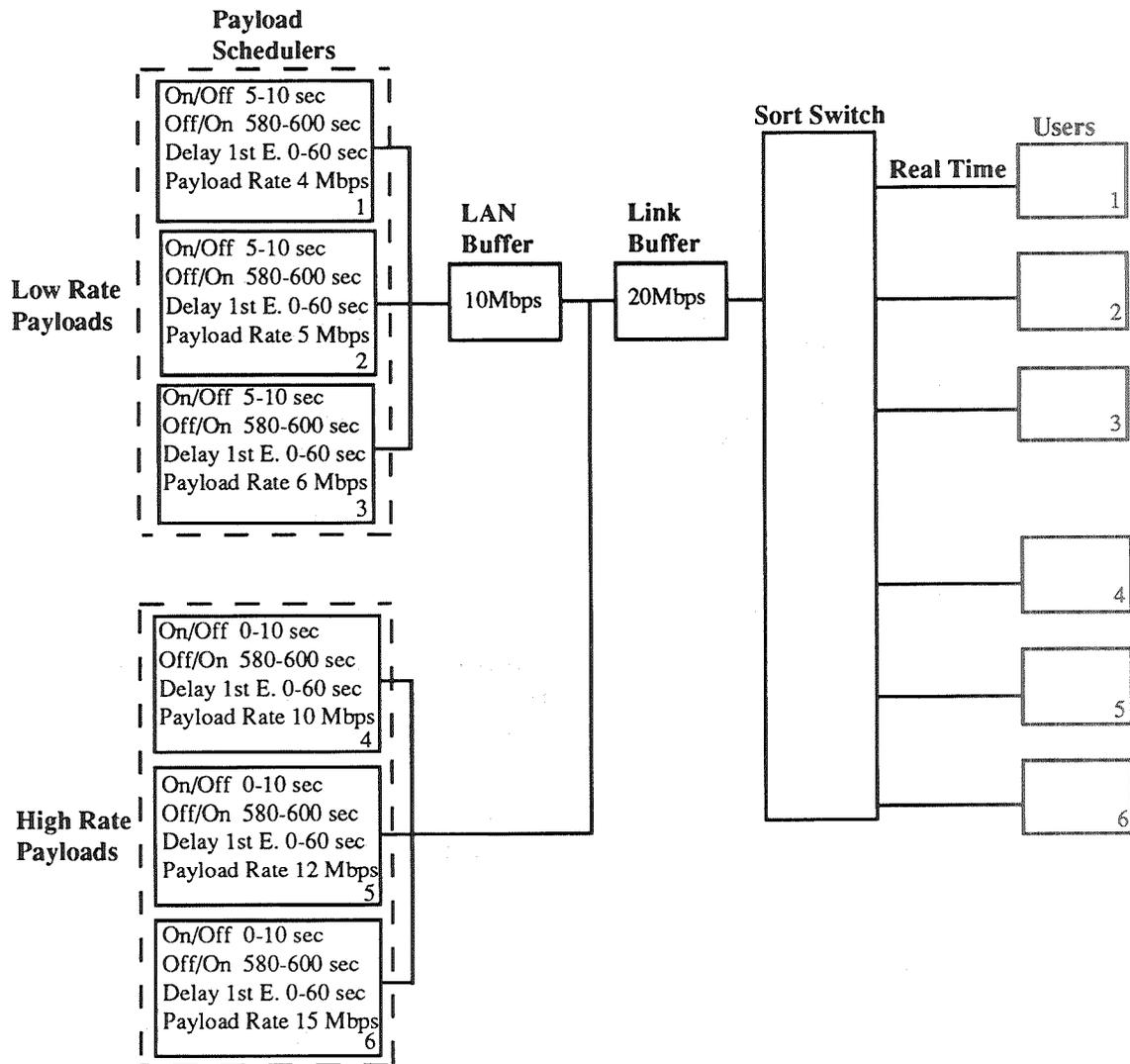
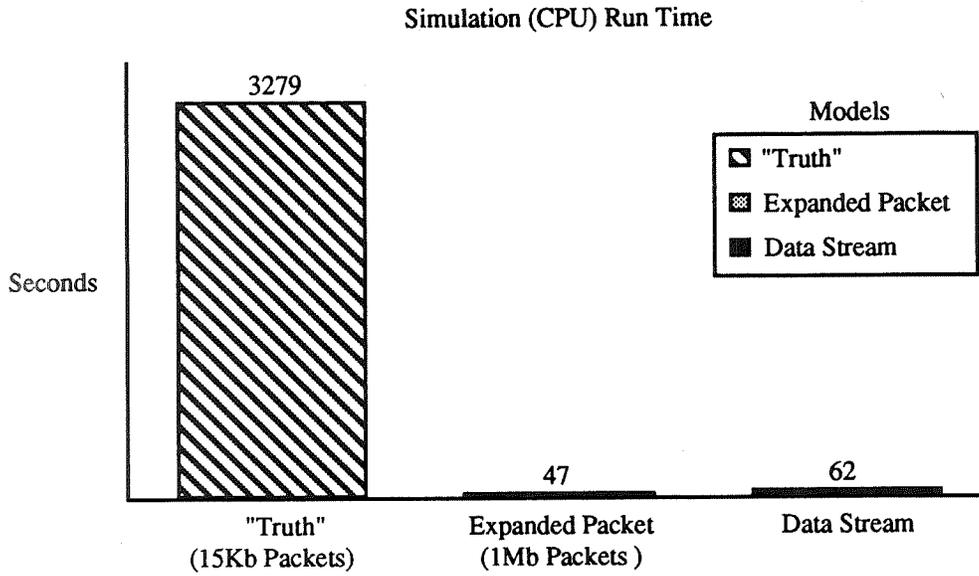


Figure 2: System used to compare data streams, artificially inflated packet and truth models.

Figure 4 illustrates the percentage of difference for the data's mean transit times from source to sink of the expanded packet model and the data stream model when compared to the truth model. As can be seen, there is a significant amount of error introduced in the mean transit times of the data when the packet size is artificially increased. The mean transit time errors in the data stream model were negligible when compared to the expanded packet model for all six payloads.

A second experiment was also designed to determine the reason for the mean transit time errors found when artificially expanding the packet sizes. It was speculated that the error was due largely to the amount of queuing experienced in the system, so the second system was designed so the data reached the user not in near real time, but within an hour, which would reduce the amount of queuing in the system. The three models (truth, expanded packet, and data stream) were again constructed and executed. Figure 5 shows that the percentage of difference when comparing the expanded packet and the data stream model models to the truth model is not as significant in the one hour data delivery system as the near real time system. However, in both systems the data stream model had less error than the expanded packet model.



Model

Figure 3: Comparison of simulation times of three different models.

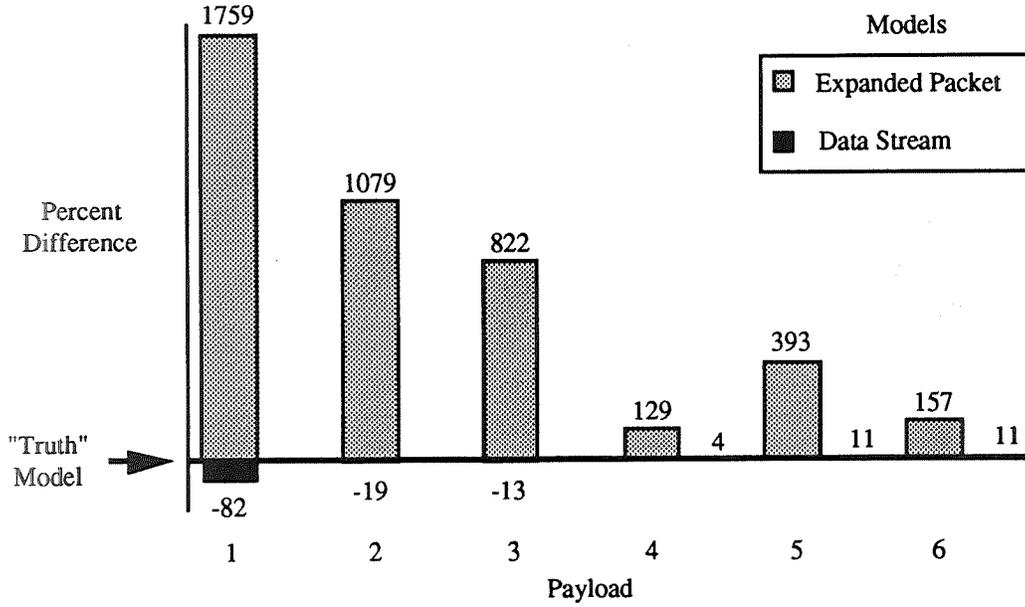


Figure 4: Mean Transit Time Difference From "Truth" Model

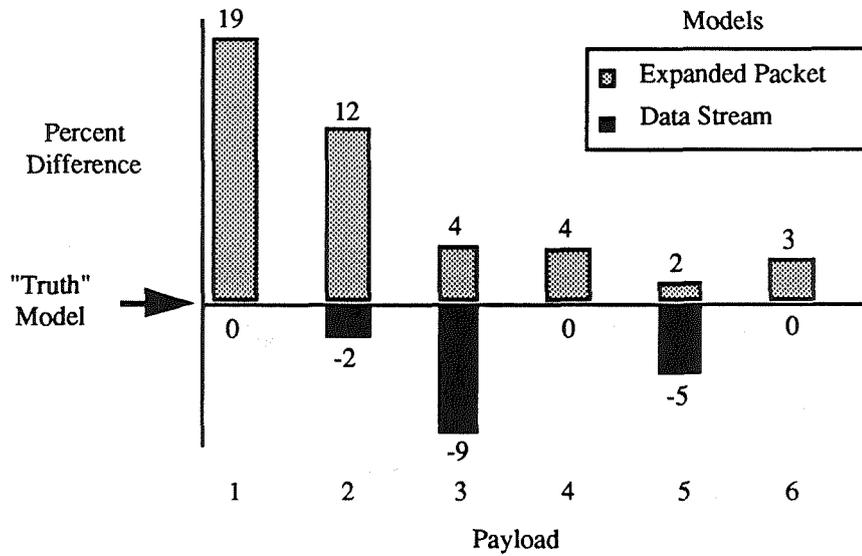


Figure 5: Mean transit time difference for one hour data delivery system.

COMMERCIAL POTENTIAL

Modeling has always been a valuable method for determining and validating requirements and designs for network-based systems. This is becoming even more important today. Competition in bids for commercial, military and other government contracts and a decrease in available make simulations of networks during the requirements and design phases of systems more important. In addition, because the systems are becoming larger and more complex, they are becoming harder to understand, and it is more difficult to predict their performance. By modeling these systems early problems can be detected. Modeling saves time and money in the long term by helping to identify potential problems (e.g. bottle necks). Lastly, detecting problems before development has the potential for large savings. The earlier errors are detected in a design, the less expensive it is to fix.

The commercial potential of DSDS is evident when compared with recent releases of other modeling tools. Other simulations use only the packet model simulation and not the data streams method. These systems require more simulation time and are more prone to error for the modeling of large end-to-end systems. Though the current commercial tools do not support the data streams methodology, vendors have expressed a desire to incorporate it into their products.

One commercial application of data streams is obvious: they can be used to simulate data management systems, whether NASA related or not. Packets can be used to represent entities in the general sense; people, cars, planes, boxes, bags or pencils to name a few. Similarly, data streams can be used to model the flow of people, cars, planes, boxes, bags and pencils. Data stream modeling is not intended to replace packet modeling. However, data streams can be used effectively as an alternative to artificially inflating the size of packets.

CURRENT WORK

Work on DSDS is still being performed. We are currently adding a graphical user interface, investigating an expanded use of data streams, and developing of an integrated simulation system. The following paragraphs discuss each of these areas.

Graphical User Interface

The current user interface for DSDS is a menu-based system where users type in the name of each element and the names of other elements that the first is connected to. A new graphical user interface is being developed that will allow users to draw the elements of the model and physically connect them together on a bitmapped screen. This will allow modelers to see their model as they develop it. The new interface will also support hierarchies so that the model can be developed in a structured, top-down or bottom-up fashion. Figure 6 shows an example of a model drawn using the new graphical user interface.

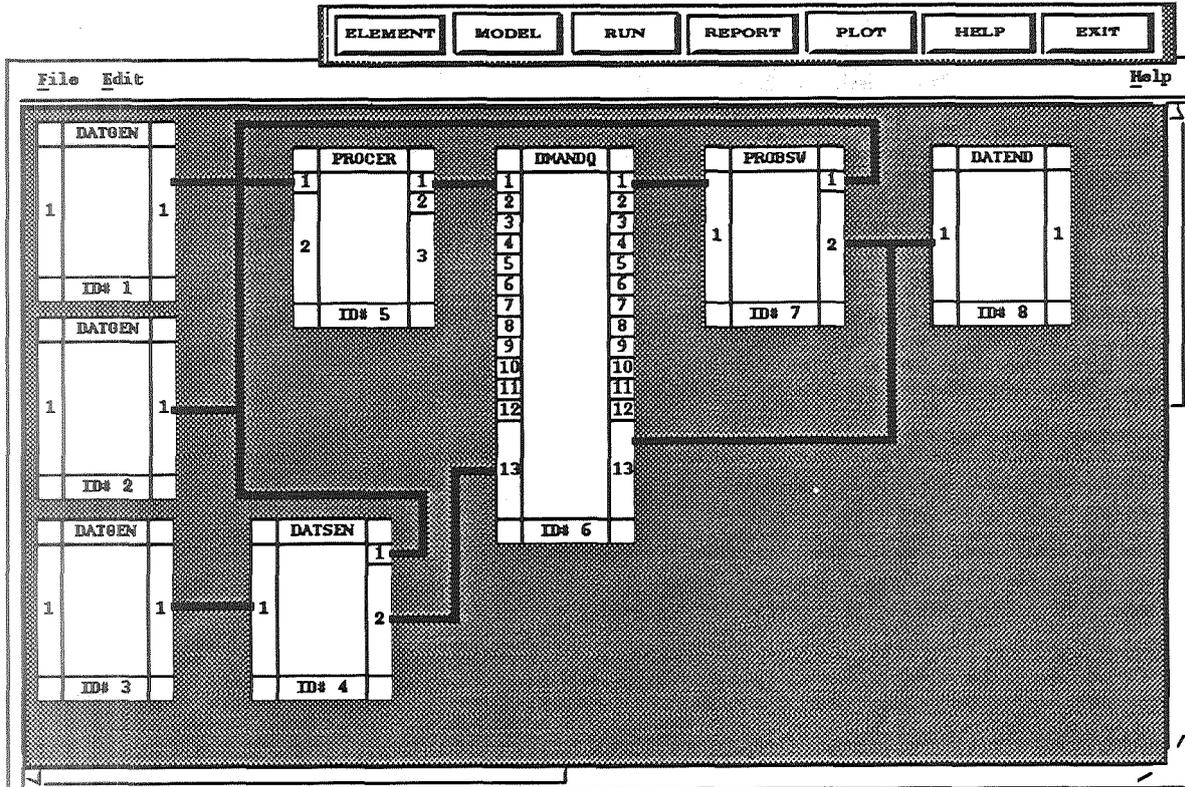


Figure 6: New graphical user interface for DSDS.

General Purpose Modeling

To date, the use of data streams has been limited to modeling of large, high-level, end-to-end data systems. These models have been characterized by duty cycles which produce millions of packets in a continuous mode. A recent study by Code 520 at GSFC has also demonstrated the potential of data streams for general purpose modeling. The study concentrated on the suitability of data streams for non-linear arrival patterns. Simulations were run with varying numbers of packet arrivals being simulated as streams. As expected, the streams which emulated the least number of arrivals gave the best results. The study reinforced the conclusion that data streams, when used for the purpose of optimizing CPU run time, are more accurate than artificially-inflated packets. In any case where the real packet size must be violated, the data stream option should be considered.

Signal Analysis Modeling System (SAMS)

An integrated modeling environment called the Signal Analysis Modeling System is also being developed to support pre-mission specification, storage, archival and analysis of spacecraft signal data. To achieve this, SAMS contains a relational database system (Oracle), an Experiment Scheduling System (ESS), a simulation system (DSDS), and a desktop publishing system (Asterix). The Oracle database contains thousands of records describing the commands to be sent to the spacecraft, together with the measurements to be collected by the onboard sensors in response to the command stimulus. This database information is exported to the scheduling system, which generates a timeline of operations by scheduling the use of the

onboard activities in such a way as to resolve conflicts. The ESS timeline is then exported to DSDS and used as a driver to populate a simulation model of the end-to-end NASA data system. The DSDS model then yields insight into the performance of the individual and aggregate data processing and data communication sub-systems. Finally, the Asterix desktop publishing system is used to capture the input data and output products generated by the other tools, such that a single, comprehensive report can be produced documenting the results of all the analyses performed.

CONCLUSION

DSDS has been used for a number of years on NASA projects but has not been released to the general public. The combination of the packet modeling, data streams modeling and extensibility makes the tool versatile and suitable for many different modeling situations. The data streams optimization technique allows accurate modeling of high data rate systems that could not be modeled accurately using the expanded packet optimization technique. As more high data rate systems are used in government and commercial applications, a technique such as data streams will become essential for modeling their performance in a timely and accurate manner.



omit

**MANUFACTURING TECHNOLOGY
PART 2**

PRECEDING PAGE BLANK NOT FILMED

394 ~~INTENTIONALLY BLANK~~



N 9 3 2 5 6 0 3

A NOVEL OPTICAL/DIGITAL PROCESSING SYSTEM FOR PATTERN RECOGNITION

Bradley G. Boone and Oodaye B. Shukla

Johns Hopkins University Applied Physics Laboratory
Electro-Optical Systems Group
Johns Hopkins Road, Laurel, MD 20723-6099

342-61

150572

P. 10

ABSTRACT

This paper describes two processing algorithms that can be implemented optically: the Radon transform and angular correlation. These two algorithms can be combined in one optical processor to extract all the basic geometric and amplitude features from objects embedded in video imagery. We show that the internal amplitude structure of objects is recovered by the Radon transform, which is a well-known result, but, in addition, we show simulation results that calculate angular correlation, a simple but unique algorithm that extracts object boundaries from suitably thresholded images from which length, width, area, aspect ratio, and orientation can be derived. In addition to circumventing scale and rotation distortions, these simulations indicate that the features derived from the angular correlation algorithm are relatively insensitive to tracking shifts and image noise. Some optical architecture concepts, including one based on micro-optical lenslet arrays, have been developed to implement these algorithms. Simulation test and evaluation using simple synthetic object data will be described, including results of a study that uses object boundaries (derivable from angular correlation) to classify simple objects using a neural network.

1. INTRODUCTION

Optics and pattern recognition are key areas for systems development for both DoD and civilian applications, including tactical missile guidance, strategic surveillance, optical parts inspection, medical imaging and non-destructive evaluation. Both passive imaging sensors (infrared (IR) and visible) and active microwave imaging sensors have been employed in many systems to date, but pattern recognition solutions in conjunction with these sensors are highly application dependent and have required extensive training. These factors have precluded extensive development.

The objective of this paper is to describe the concept of an optical processor for object measurement that can be interfaced to a variety of sensors, including imaging IR, optical machine vision systems and synthetic aperture radar (SAR), thus making it very versatile. The optical processor will be used in conjunction with a neural network algorithm to classify objects. Another goal is to provide a preliminary report on the effectiveness of the measurement and classification algorithms that we plan to implement with optics and digital electronics, respectively.

One of the key concerns in the design of optical feature extractors and image matchers for target recognition is that performance be invariant with respect to position, scale and rotation distortions. In many cases traditional approaches have involved complicated mathematical transformations to achieve distortion invariance¹⁻². Our approach³⁻⁵ offers a compact distortion-insensitive method of optical correlation using primitive image operations such as image replication, multiplication, integration, and detection, and is useful in viewing objects in plan-view. One of the key aspects of this approach is that we use optical correlation to measure objects rather than match them. We leave the matching up to a neural network.

Our optical processor⁵ concept implements the optical Radon (or Hough) transform and the APL developed optical angular correlation technique, followed by appropriate numerical processing, and a neural net classifier. The angular correlator is a unique development that enables object symmetry, orientation, primitive dimensions and boundary to be estimated. Along with the well-developed Hough transform, which provides information on the internal structure of objects, these features collectively describe most simple closed-boundary objects in an elegant and compact way, thus affording generic object measurement and the prospect of effective object classification. The only major requirements for this processor are that the input imaging sensor employ

396

INTENTIONALLY BLANK

detection with adaptive thresholding and centroid tracking, both of which are common (or easily implemented) attributes of most imaging sensors.

The key components of this overall concept are laid out in Figure 1. Basically, the overall system consists of four stages: an optical interface to an appropriate sensor display or entrance optics, optical processor for angular correlation and (optional) Hough transform, digital processor for calculating various features of the object data, and finally, a digital or analog neural network. Either of two techniques have been considered for rotating input imagery: video feedback (time-multiplexing) or multiple aperture optics (space-multiplexing). (However, image rotation is not critical to implementing the angular correlation algorithm.) Early efforts⁶⁻⁷ using video feedback explored the Hough transform and angular correlation in a single-channel (time-multiplexed) implementation. Similar concepts in multi-aperture optical processors have been proposed for these and other applications⁹⁻¹² and offer smaller size and increased throughput.

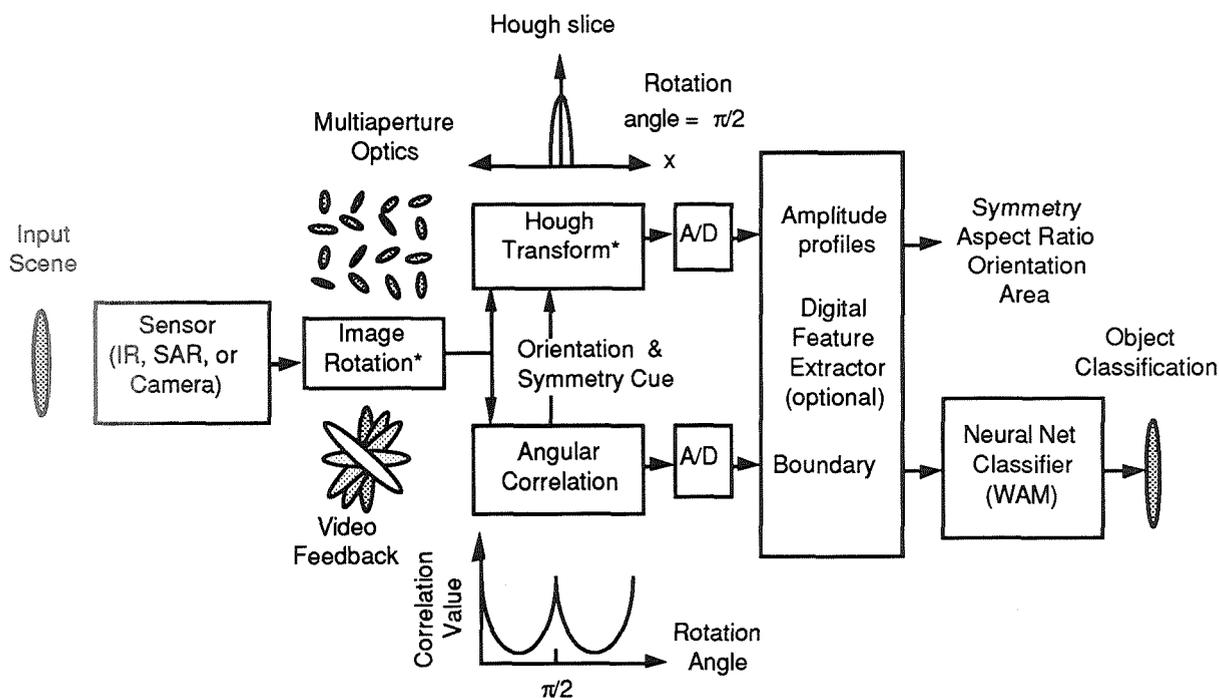


Figure 1. Block diagram of angular correlator architecture (items denoted by * are optional).

What we will concentrate on in this paper is the evaluation of a simulation of the angular correlation algorithm in conjunction with a winner-take-all associative memory (WAM) neural net. Both of these algorithms were developed at JHU/APL and are currently in the process of implementation. Some aspects of the implementation of the angular correlator will also be described.

2. ANGULAR CORRELATION

Many kinds of correlation algorithms have been implemented for pattern matching applications. Most correlation algorithms shift one image with respect to the other while calculating the area of overlap. Angular correlation simply calculates the area of overlap versus angle. The resulting set of correlation values can be used to recover the boundary of an object if the object is thresholded and binarized and the other object is a slit. Since the angular correlation algorithm uses rotation, the need for rotational invariance is obviated. Scale invariance is irrelevant to angular correlation because it measures an object. The optical implementation uses incoherent light (allowing the use of standard camera optics or video displays as image inputs) and multi-aperture optics (making it light-weight and compact using off-the-shelf components).

Primitive features of an object can also be determined by using angular correlation. Object primitives include: area, length, width, aspect ratio, symmetry, and orientation. For most simple objects periodicity of the boundary is directly related to symmetry. For a square (which has four-fold symmetry), the periodicity is $\pi/2$ ($2\pi/4$), whereas for an equilateral triangle (three-fold symmetry), the periodicity is $2\pi/3$. With the object centered, the peak values of the recovered boundary are related to the maximum extent of the object. The minimum value of the boundary curve is the minimum dimension of the object through its centroid. Taking the ratio of the maximum and minimum values of the recovered boundary yields the aspect ratio of a simple two-fold symmetric object like a rectangle or an ellipse. The minimum value of the recovered boundary curve (or "bias") is also a measure of the image "mass" concentration of the object about its centroid. For example, the boundary of a star shaped object is a set of periodic peaks with a lower bias than the boundary of a square which has more image "mass" at its center (see Figure 5).

For angular correlation, both objects have to be centered with respect to a common origin. In that case the maximum correlation value gives the cue for selection of an optimal Radon transform slice. The peak correlation values for offset slits are less than the peak correlation value of a slit with no offset as shown in Fig. 2(a). Even for large offsets the periodicity and optimum cueing angle remain unchanged. Thus a key assumption necessary for implementing this algorithm digitally or optically is that the object be centroid tracked, something often achieved in practice with imaging sensors and good tracking systems.

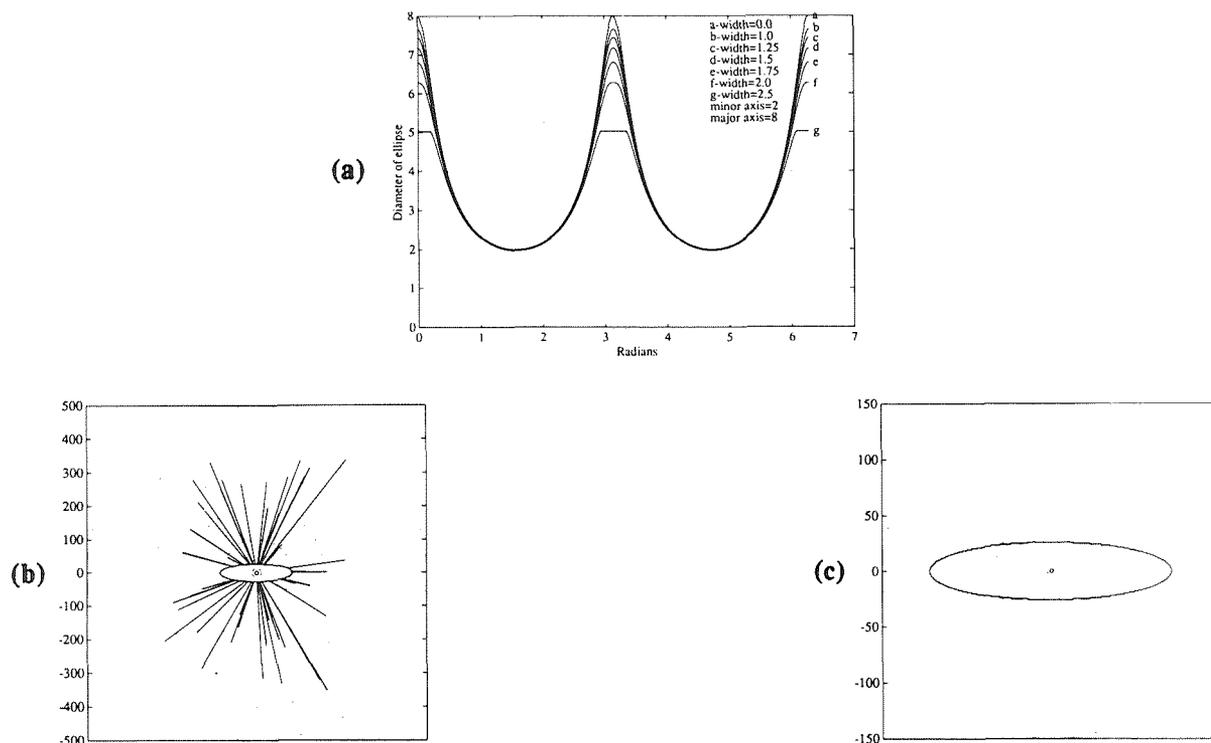


Figure 2. (a) Angular correlation of ellipse for varying offsets of correlation slit, (b) polar plot of ellipse with substantial image noise ($SNR \approx 8$), and (c) cleaned-up ellipse boundary using a thresholded first-forward difference filter.

Consider the angular cross-correlation of a slit (extending over the entire image) with an ellipse in more detail. When the slit and the major axis of the ellipse are oriented horizontally, the area of overlap is approximately the width of the slit multiplied by the major axis of the ellipse. By decreasing the slit width, the intersecting area approaches the length of the major axis of the ellipse. As the slit is rotated, the area of overlap is the boundary of the object at the slit angle. Angular cross-correlation for binarized images can then be expressed mathematically as:

$$R(\theta) = \int_0^{\infty} \int_0^{\infty} \text{rect}\left[\frac{r}{r(\theta')}\right] \text{rect}\left[\frac{r}{r_{\text{rect}}(\theta' + \theta)}\right] r dr d\theta' \quad (1)$$

where $\text{rect}[r/r_{\text{rect}}(\theta' + \theta)]$ is a functional description of a rectangular slit rotated by an angle θ , and $\text{rect}[r/r(\theta')]$ is the corresponding description for the desired object.

Ideally the slit width should approach zero to recover the exact boundary of an object, but in practice we must measure a finite signal. The minimum sampling angle necessary to sample the boundary of an object and to satisfy the Nyquist sampling criterion can be calculated by examining the Fourier spectrum of the boundary function to obtain the cutoff frequency. Then the appropriate slit width can be determined from this using simple geometry. Essentially this means that complex binarized objects should be cross-correlated with a slit one pixel wide. In the digital simulation, the slit is rotated while the image remains fixed. If the image is noise-free, then the boundary recovered is exact. However if the image is extremely noisy, then the recovered boundary will have spikes on it as shown in Figure 2(b). A noisy boundary can be filtered to recover the smooth boundary by using the Nyquist bandwidth of the boundary function to set a low-pass filter cutoff (or by using a first-forward difference with a limiting threshold as was done for Fig. 2(c)).

The angular correlation algorithm is effective on simple convex shapes, such as rectangles, triangles, ellipses, and circles, and some concave shapes such as stars and gears. Objects that do not have simple closed boundaries are those with re-entrant boundaries and multiple boundaries, as shown in Figure 3. For objects with such boundaries, the estimated boundary recovered by angular correlation will not necessarily enable it to be discriminated from other (simpler) boundaries because only the total area of overlap is recovered. In other words, the area of overlap between the slit and object is a single value that does not contain any information about boundaries within the slit.

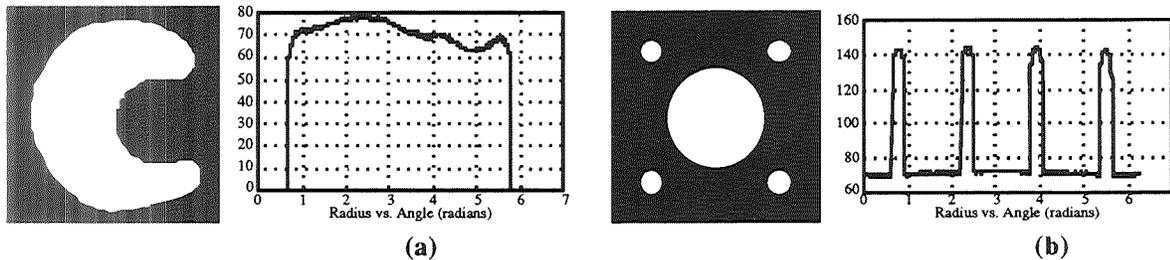


Figure 3. Objects and boundary functions of shapes with (a) re-entrant boundaries and (b) multiple boundaries.

3. HOUGH TRANSFORM

The Hough (or Radon) transform is a well known mathematical transform used in image processing to reconstruct objects. The Radon transform is a collection of 1-D projections. For each angle θ , an object's amplitude projection is obtained by integration perpendicular to the p-axis (which is the x-axis rotated by θ). The complete Radon transform is given by:

$$F_R(p, \theta) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(r) \delta(p - \mathbf{r} \cdot \mathbf{n}) d^2r \quad (2)$$

where the 2-D object is defined by the function $f(r)$ and \mathbf{n} is the unit vector normal to the p-axis. The (p, θ) coordinates represent Radon space. An interesting connection can be drawn between the Radon transform and correlation. If, instead of calculating the angular correlation we calculate the annular correlation. i.e.:

$$R(r) = \int_0^{r_m} \int_0^{2\pi} f(r, \theta) \text{rect}\left(\frac{r' + r}{r_0}\right) r' dr' d\theta \quad (3)$$

we can obtain a result that is equivalent to the Radon transform averaged over all θ . For simple 2-fold symmetric objects (like a rectangle or ellipse) embedded in backgrounds that can be well-thresholded the result is very nearly the same as a single optimal Radon slice. This result is particularly useful for recovering the re-entrant and multiple-boundary objects mentioned earlier. The combined annular correlation (Radon transform) and angular correlation yield the correct reconstructions for these two cases as shown in Fig's. 4(a) and (b).

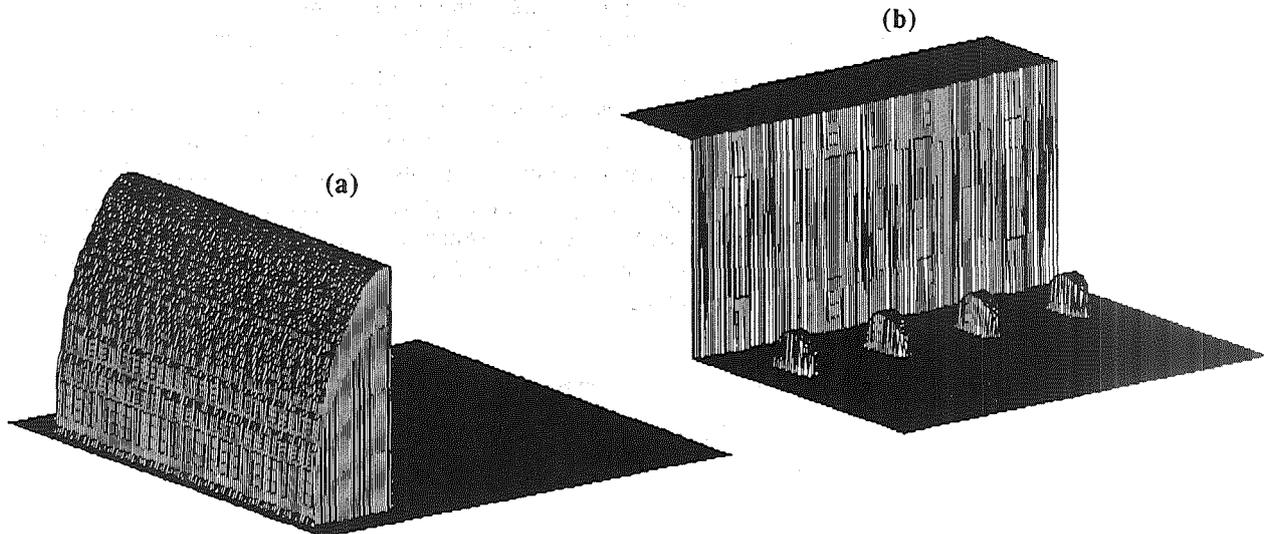


Figure 4. (a) Reconstructed re-entrant boundary shape and (b) multiple boundary shape resulting from the use of angular and annular correlation (or Radon transform).

The Hough transform can also be used to classify objects compared with or without angular correlation, especially 2-fold symmetric objects. If angular correlation is used the angle corresponding to maximum correlation is used to determine the orientation of the object containing the maximum information from the Hough transform. The internal amplitude profile of the object along this orientation is usually the optimal Hough transform slice. Given the orientation and the corresponding profile of an object (or its Fourier components), classification algorithms have been used to identify the object³. Primitive features of the object such as length, width, and aspect ratio and orientation derived from angular correlation have also been used along with the Hough transform to improve neural net classification rates³.

4. SIMULATION OF ANGULAR CORRELATOR

Before implementing the angular correlator optically, it is useful to simulate it on a computer using a set of simple synthetic geometric objects shown in Figure 5. These objects were chosen to include four-fold, three-fold, and two-fold objects as well as an object with multiple boundaries (Figure 5(d)) and an object with a re-entrant boundary (Figure 5(p)). For objects with multiple boundaries, the recovered boundary curve shows the outline of all the objects but not the boundaries between distinct constituent objects. The four-fold objects were chosen to compare shapes with differing image "mass" concentrations (square, plus-sign, rotor, and star) and perturbations from a square (trapezoid, parallelogram, and quadrilateral). The three-fold objects were arranged to compare on the basis of image "mass" concentrations (triangle and three-pointed star) and evidence of bilateral symmetry (isosceles triangle and triangle with concave sides). The two-fold objects were chosen to reflect differences in boundary frequency content (rectangle versus ellipse) and phase (ellipse versus rotated ellipse). For the object with the re-

entrant boundary, the angular correlation algorithm recovers a straight line approximation of the interior concave sides (Figures 3(a) & 5(p)). For the multiple boundary object (five circles), the correlation algorithm does not detect disjoint boundaries (Figures 3(b) & 5(d)). In these two cases, the recovered boundary from the angular correlation algorithm is not sufficient to calculate the exact primitive features of the object or objects within the image. However, as stated before, the Hough transform may be used to recover the internal structure of these objects. Otherwise boundaries of the remaining objects are recovered successfully.

5. OPTICAL HARDWARE

The angular correlator and the Hough transformer can be implemented optically in two basic architectures: time multiplexing (video feedback) or space multiplexing (multiple lenslet arrays). Although the feedback approach was first implemented as described previously³⁻⁴, the preferred implementation is a multi-aperture micro-optical architecture. A multi-aperture optical system to optically rotate an image, calculate its Hough transform and recover its boundary using angular correlation has been conceived⁵, and the experimental breadboard is shown in Figure 6. This breadboard includes a video display and collimating lens that serve as the optical interface to represent object space. The actual optical processor is preceded by a zoom lens and (optional) microchannel plate. The microchannel plate forms a real image to be replicated by the multi-lenslet array. It can also be used in a saturated mode to binarize the object. Alternatively the original video display can project through (as a virtual object), or a fiber optic window or binarizing spatial light modulator can be used to create a displayed image. The replicated images are passed through a fixed mask onto a multiple detector array as shown in Fig. 6 (inset(a)). Each detector spatially integrates the superposition of each replicated image and the

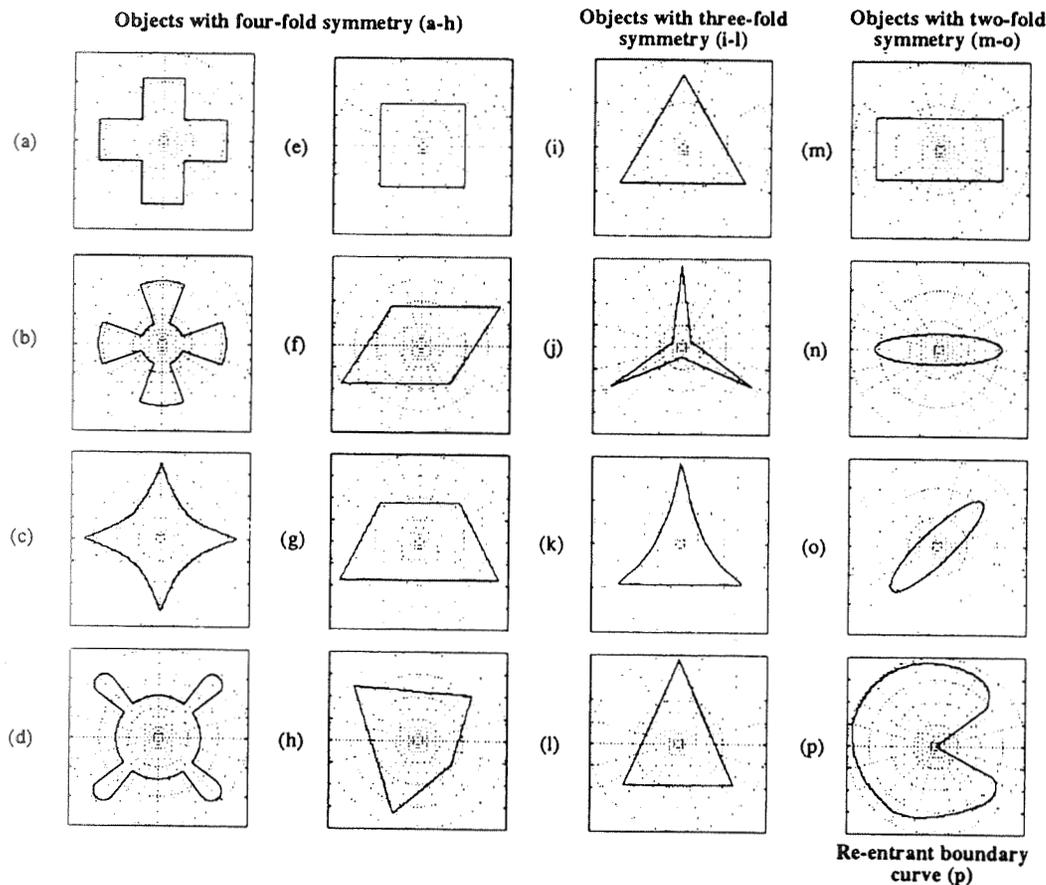


Figure 5. Object boundaries recovered by angular correlation, grouped by symmetry and used to evaluate the neural network simulation.

corresponding mask pattern. For angular correlation the mask consists of a series of rotated half-plane slits as shown in Fig. 6 (inset(b)). For annular correlation (θ -averaged Hough transform) it is a series of annuli as shown in Fig. 6 (inset(c)). Ordinarily, in order to implement the Hough transform, the image has to be rotated. Previous optical architectures use mechanical rotation schemes⁶, but mechanical devices have reliability problems and will always be a throughput bottleneck. Replicating the image optically and rotating either the images optically or the pattern elements of the processor mask will substantially increase the throughput and reliability. Other functions can also be performed, such as tracking, using simple mask patterns. The detector outputs are then preamplified, filtered, multiplexed and A/D converted to input into a PC-hosted neural network algorithm discussed in the next section.

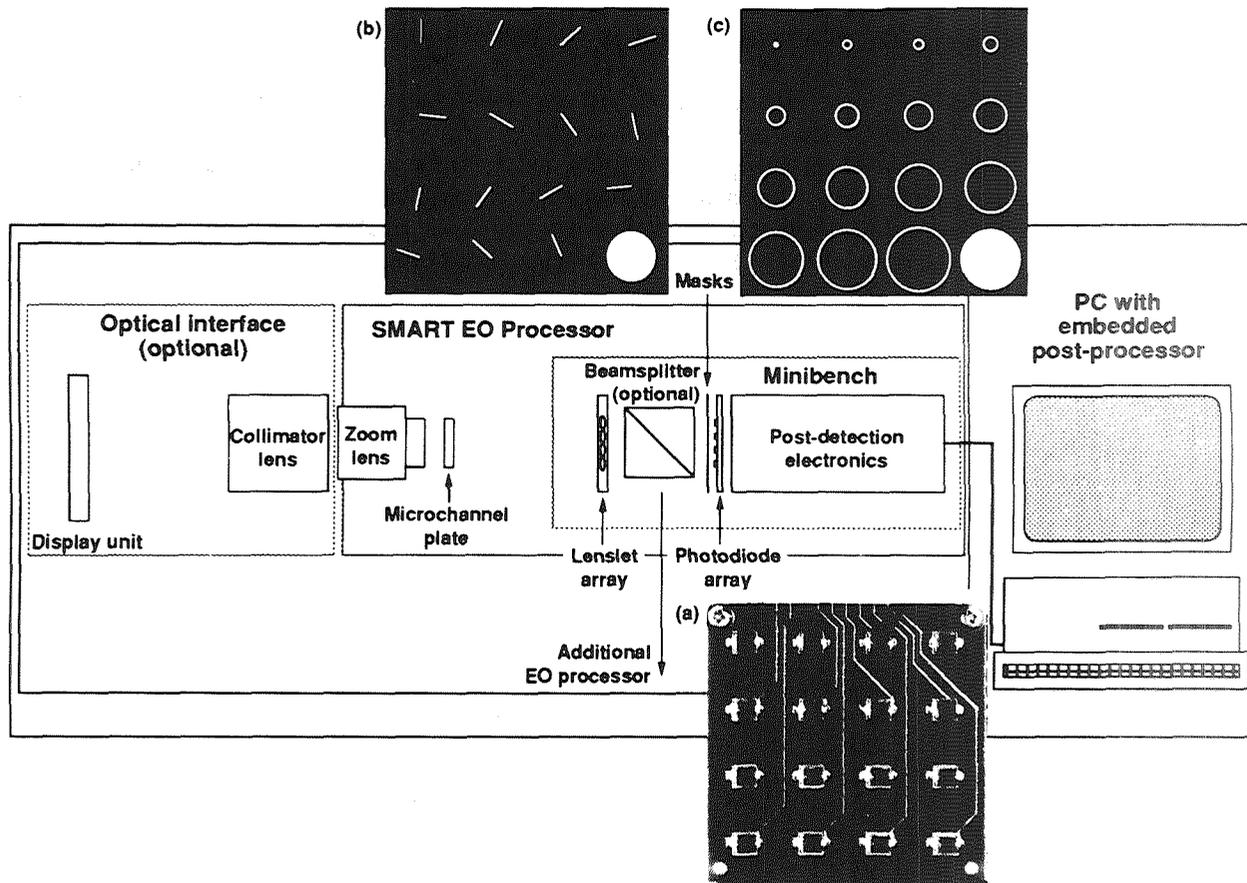


Fig. 6. Optical hardware layout for experimental breadboard optical/digital processor for pattern recognition, showing detector plane (inset a) and masks for angular and annular correlation (insets b and c)

6. DIGITAL NEURAL NETWORK CLASSIFIER EVALUATION

We can easily distinguish between two-fold, three-fold, and four-fold symmetric objects, i.e.: differentiate a rectangle, triangle, and square using symmetry alone. The next higher level of discrimination within the same symmetry class is "bias" or the minimum value of the boundary curve. At this level, image "mass" concentrations (spatially integrated intensities) of objects are compared, as given by the ratio of minimum to maximum radii of the object. The lower the ratio, the greater the image "mass" concentration. (with a minimum of 1 for the circle). For instance, at this level we can separate a rotor from a star. Although it won't be discussed here, Fourier spectra of the boundary can also be used to categorize these objects. This approach is similar to Fourier boundary descriptors described elsewhere¹³

described elsewhere¹³

In this section we will focus on a neural network that uses the boundaries obtained by angular correlation for classification. The angular correlation outputs will be matched against patterns stored in a digital simulation of an analog VLSI neural network under development at APL. This neural network paradigm is a modification of the Bi-directional Associative Memory (BAM) developed by Kosko¹⁴ in which the hidden layer is replaced by a winner-take-all layer¹⁵. Using the winner-take-all layer eliminates the need for the output to be fed back to the input and eliminates convergence problems prevalent in feedback networks. In the winner-take-all layer, one neuron dominates all other neurons within the layer and allows only one match to occur, thus removing the ambiguity in pattern matching applications. The input patterns are stored as weights of the hidden layer in a bipolar format (off-bits represented by -1 and on-bits represented by +1). The input pattern is fed into the winner-take-all associative memory (WAM), and a binary code representing the matched pattern is the output.

The operation of the WAM can be explained using matrix multiplication. The boundary patterns generated by angular correlation are thresholded, binarized and stored in a pattern matrix. To match the stored patterns against an input pattern, the inner product of the transpose of input vector and the pattern matrix is calculated. An input boundary can then be compared to all of the stored boundaries by forming the correlation matrix. The resulting correlation matrix is a measure of how similar individual stored patterns are to other stored patterns in the pattern matrix. If the correlation matrix has all unity diagonal terms and zero off-diagonal terms, then each input pattern matches itself exactly and is perfectly distinguishable from the other patterns within the pattern matrix. For more details see reference (16).

Figure 5 shows the objects used to test this network. The objects were first presorted according to their symmetry. The re-entrant curve was placed in its own class and was compared with all the objects in the other three symmetry classes. The boundary curves of the presorted objects were binarized with three different threshold levels: median, high and low and stored as rows in separate pattern matrices within each symmetry class. The patterns are also subsampled to 40 data points to simulate more closely the actual VLSI implementation of this network, which stores 116 patterns each 124 bits long. Each of the thresholded versions of each boundary pattern were concatenated, padded with zeros and fed in as a pattern of 124 bits into the WAM simulation. The resulting correlation matrices for the four-fold (A), three-fold (B), two-fold (C) classes, and the re-entrant (D) class are:

$$\mathbf{A} = \begin{matrix} & \begin{matrix} \text{(a)} & \text{(b)} & \text{(c)} & \text{(d)} & \text{(e)} & \text{(f)} & \text{(g)} & \text{(h)} \end{matrix} \\ \begin{matrix} \text{(a)} \\ \text{(b)} \\ \text{(c)} \\ \text{(d)} \\ \text{(e)} \\ \text{(f)} \\ \text{(g)} \\ \text{(h)} \end{matrix} & \begin{bmatrix} 1.0000 & 0.5167 & 0.5833 & 0.7333 & 0.7833 & 0.6167 & 0.6167 & 0.5667 \\ 0.5167 & 1.0000 & 0.4667 & 0.3167 & 0.3000 & 0.4333 & 0.4333 & 0.4500 \\ 0.5833 & 0.4667 & 1.0000 & 0.4500 & 0.3667 & 0.9667 & 0.9667 & 0.9833 \\ 0.7333 & 0.3167 & 0.4500 & 1.0000 & 0.9167 & 0.4500 & 0.4500 & 0.4667 \\ 0.7833 & 0.3000 & 0.3667 & 0.9167 & 1.0000 & 0.4000 & 0.4000 & 0.3833 \\ 0.6167 & 0.4333 & 0.9667 & 0.4500 & 0.4000 & 1.0000 & 0.9667 & 0.9500 \\ 0.6167 & 0.4333 & 0.9667 & 0.4500 & 0.4000 & 0.9667 & 1.0000 & 0.9500 \\ 0.5667 & 0.4500 & 0.9833 & 0.4667 & 0.3833 & 0.9500 & 0.9500 & 1.0000 \end{bmatrix} \end{matrix} \quad (4)$$

$$\mathbf{B} = \begin{matrix} & \begin{matrix} \text{(i)} & \text{(j)} & \text{(k)} & \text{(l)} \end{matrix} \\ \begin{matrix} \text{(i)} \\ \text{(j)} \\ \text{(k)} \\ \text{(l)} \end{matrix} & \begin{bmatrix} 1.0000 & -0.1167 & 0.9000 & 0.4500 \\ -0.1167 & 1.0000 & -0.1167 & 0.3667 \\ 0.9000 & -0.1167 & 1.0000 & 0.3833 \\ 0.4500 & 0.3667 & 0.3833 & 1.0000 \end{bmatrix} \end{matrix} \quad (5)$$

$$\mathbf{C} = \begin{matrix} & \begin{matrix} \text{(m)} & \text{(n)} & \text{(o)} \end{matrix} \\ \begin{matrix} \text{(m)} \\ \text{(n)} \\ \text{(o)} \end{matrix} & \begin{bmatrix} 1.0000 & 0.3167 & 0.0333 \\ 0.3167 & 1.0000 & 0.0833 \\ 0.0333 & 0.0833 & 1.0000 \end{bmatrix} \end{matrix} \quad (6)$$

$$\mathbf{D} = \begin{matrix} & \begin{matrix} \text{(a)} & \text{(b)} & \text{(c)} & \text{(d)} & \text{(e)} & \text{(f)} & \text{(g)} & \text{(h)} & \text{(i)} & \text{(j)} & \text{(k)} & \text{(l)} & \text{(m)} & \text{(n)} & \text{(o)} \end{matrix} \\ \begin{matrix} \text{(a)} \\ \text{(b)} \\ \text{(c)} \\ \text{(d)} \\ \text{(e)} \\ \text{(f)} \\ \text{(g)} \\ \text{(h)} \\ \text{(i)} \\ \text{(j)} \\ \text{(k)} \\ \text{(l)} \\ \text{(m)} \\ \text{(n)} \\ \text{(o)} \end{matrix} & \begin{bmatrix} 0.667 & 0.214 & 0.262 & 1 & 0.929 & 0.262 & 0.262 & 0.286 & 0.357 & -0.476 & 0.262 & 0.190 & 0.333 & -0.548 & -0.190 \end{bmatrix} \end{matrix} \quad (7)$$

Under these conditions the patterns are clearly distinguishable from each other. For each column, there is a single maximum value. These results show that if multiple thresholds are used, then boundary curves can be sub-sampled and still be matched using the WAM network. The number of samples needed to adequately sample the boundary function determines the number of channels required for the multi-aperture optical architecture. Since the complexity of the multi-aperture design decreases as the number of channels decrease, sub-sampling the boundary functions and using three thresholds minimizes the number of optical channels needed to compute the angular correlation. Having fewer optical channels also simplifies the interface electronics needed to transfer the detector signals to the neural network. Since the storage capacity of the WAM is limited, as the number of thresholds increases, the amount of boundary samples must decrease. The optimal number of thresholds and the minimum number of samples required is a trade-off issue that is application dependent and requires further study.

7. APPLICATIONS

Some of the applications well suited for this classifier based on angular correlation algorithm include optical parts inspection, machine vision, non-destructive evaluation, medical imaging, strategic surveillance, and tactical missile guidance. For missile guidance, this classifier can be employed in a SAR sensor to recognize targets such as ships. In strategic surveillance the angular correlator can be mated to one of several different sensors which include optical, infrared, SAR, or microwave radiometric sensors. For non-destructive evaluation, the patterns contained in the surface morphology of the test object can be measured using the angular correlator and primitives derived from it. In the area of medical imaging the growth, shape and size of retinal lesions can be measured and tracked over time using this algorithm to recover the boundary of the diseased area. In machine vision, a robot can use this algorithm to recognize, avoid, or manipulate objects using boundaries.

For optical parts inspection, parts traveling down a conveyor belt can be imaged and replicated through the optical interface to extract their boundaries and then sent to the WAM. The WAM matches the part's boundary to one of the stored boundaries. Further action may be taken depending on whether there is a match or not, such as sorting. The parts inspection or sorting task is a well-constrained task. The placement and orientation of the moving parts can be well defined and the boundaries of the parts are well-known. The boundaries can be thresholded to obtain the best match and stored in the pattern matrix. The simplicity of this design makes this correlator compact, reliable and suitable for the factory environment. In this correlator, the throughput is limited by the electronics used to scan the detector array and the WAM. The detector array can be scanned on the order of microseconds and the processing time of the WAM is on the order of 100 microseconds. Thus the correlator should be able to classify objects on a conveyor within 100-200 microseconds, thus making it suitable for fast conveyors.

8. SUMMARY AND CONCLUSION

In this paper we have described angular correlation algorithm that can be used to obtain the boundary of an object. From the boundary, the object can be recognized using a neural network. A winner-take-all associative memory (WAM) network was used to match the object boundaries with stored boundaries. The results from the WAM simulation showed that it is possible to recognize objects based on their boundary. The angular correlation algorithm can be easily implemented using multi-aperture optics to replicate the input object and cross-correlate it with a series of rotated slits. This implementation is well-suited to interface with the WAM for classification. With multiple thresholds, the boundary curves can be decimated and still be recognized by the WAM. Even modified or perturbed objects can be identified using multiple thresholds. The boundary of the object is derived optically in parallel almost instantaneously, and the throughput is limited to less than 100-200 microseconds by the read-out electronics and the WAM classifier. One of the many applications for this hybrid processor is optical parts inspection where a set of objects with known boundaries can be matched to the stored boundaries.

9. ACKNOWLEDGMENT

Discussions with R. E. Jenkins on the neural network simulation and support from J.R. Connelly and W.S. Denny on designing and fabricating detector electronics are gratefully acknowledged.

10. REFERENCES

1. D. Casasent, and D. Psaltis, "New Optical Transforms for Pattern Recognition", Proc. IEEE 65 , 77 (1977)
2. D. Casasent, and D. Psaltis, "Position, rotation and scale invariant optical correlation", Applied Optics 15, 1795 (1976)
3. B. G. Boone, O. B. Shukla, and D. H. Terry, "Extraction of features from images using video feedback", Automatic Object Recognition, Proc. SPIE Vol. 1471, 390 (1991)
4. B. G. Boone, O. B. Shukla, and M. D. Bulla, "Method and Apparatus for Radon Transformation and Angular Correlation in Optical Processors", U. S. patent # 5,101,270, May 7, 1992
5. B. G. Boone, and O. B. Shukla, "SMART Electro-Optical Processor", JHU/APL invention disclosure, Nov. 4, 1991
6. J. P. Crutchfield, "Space-Time Dynamics in Video Feedback", Physics 10D, 229 (1984)
7. J. Cederquist, and S. H. Lee, "The Use of Feedback in Optical Information Processing", Appl. Phys, 18, 311 (1979)
8. G. R. Gindi, and A. F. Gmitro, "Optical feature extraction via the Radon transform", Opt. Eng. 23, 499 (1984)
9. K. Bromley, A. C. H. Louie, R. D. Martin, J. J. Symanski, T. E. Keenan, and M. A. Monahan, "Electro-optical signal processing module", SPIE Vol. 180, Real-Time Signal Processing II, pp. 107-113 (1979)
10. I. Glaser, "Noncoherent optical processor for discrete two-dimensional linear transformations", Optics Letters 5, 449 (1980)
11. M. Agu, A. Akiba, and S. Kamemaru, "Multimatched filtering system as a model of biological visual systems", SPIE Vol. 1014, Micro-Optics, pp. 144-150 (1988)
12. W. A. Christen-Barry, D. H. Terry, and B. G. Boone, "Detection of DNA sequence symmetries using parallel micro-optical devices", SPIE Vol. 1564, Optical Information Processing Systems and Architectures III, pp. 177-188 (1991)
13. E. L. Brill, "Character Recognition via Fourier Descriptors", WESCON, Paper 25/3, Los Angeles, CA (1968)
14. B. Kosko, "Bidirectional associative memories", IEEE Trans. Syst. Man. Cybern., Vol. 18, pp. 49-60, Jan./Feb., 1988
15. K. A. Boahen, P. O. Pouliquen, A. G. Andreou, R. E. Jenkins, "A Heteroassociative Memory Using Current-Mode MOS Analog VLSI Circuits", IEEE Transactions on Circuits and Systems, Vol. 36, No. 5, May, 1989.
16. O.B. Shukla, and B.G. Boone, "Optical Feature Extraction Using the Radon Transform and Angular Correlation", Proc. SPIE Conference on "Optical Information Processing Systems and Architectures IV", Vol. 1771, to be published, presented July 19-24, San Diego, CA

VISION-AIDED MONITORING & CONTROL OF THERMAL SPRAY, SPRAY FORMING, AND WELDING PROCESSES

John E. Agapakis
Automatix Inc.
Billerica, MA 01821

&

Jon Bolstad
Control Vision Inc.
Idaho Falls, ID 83401

N 9 3 - 2 5 6 0 4

150513
p-8

ABSTRACT

Vision is one of the most powerful forms of non-contact sensing for monitoring and control of manufacturing processes. However, processes involving an arc plasma or flame such as welding or thermal spraying pose particularly challenging problems to conventional vision sensing and processing techniques. The arc or plasma is not typically limited to a single spectral region and thus cannot be easily filtered out optically. This paper presents an innovative vision sensing system that uses intense stroboscopic illumination to overpower the arc light and produce a video image that is free of arc light or glare and dedicated image processing and analysis schemes that can enhance the video images or extract features of interest and produce quantitative process measures which can be used for process monitoring and control. Results of two SBIR programs sponsored by NASA and DoE and focusing on the application of this innovative vision sensing and processing technology to thermal spraying and welding process monitoring and control are discussed.

INTRODUCTION

Thermal spray processes, electric arc welding, laser welding, and other energy-intensive high luminosity industrial processes are normally quite difficult to monitor with the human eye because the luminous volume of the plasma or flame obscures the details of the behavior of the solid or molten material in the heat affected area. Furthermore, when one attempts to use a photographic or video camera, the viewing is further degraded by the extreme brightness variation across the image area, making it impossible to achieve proper exposure throughout the image — except possibly for small areas of comparable brightness. Optical filtering with neutral density filters — like the ones used in a welder's helmet — do not particularly help either. In the case of arc welding, one can expect to see a bright fireball at the center of the welding pool, but most of the detail at the edge of the welding pool and in the area of the welding seam and groove will be lost in relative darkness. With thermal spray processes, the injection and flow of particles within the plasma flame is almost totally concealed by the extreme brightness of the flame. In addition, the particles quickly accelerate to very high speeds, making their detection even more difficult.

Over the last six years, Control Vision has developed a unique viewing system capable of overcoming the extreme variation in scene brightness created by high luminosity phenomena such as flames, arcs, or plasmas and electronically producing a video image virtually free of arc glare. This patented system incorporates external illumination in the form of intense pulsed laser light. The laser light reflected from the site is for an instant much brighter than either the direct or reflected light of the process. The system exploits this situation by viewing the process with a special-purpose video camera equipped with a CCD video sensor and a very high speed electronic shutter synchronized with the laser flash and the framing of the video sensor.

This innovative viewing system has already been used in a variety of applications in order to provide visual feedback to an operator for in-process monitoring during production or to allow observation of phenomena not otherwise visible during process research and development. In order to incorporate such a sensor in automated process monitoring or control applications, it is necessary to develop approaches for processing and analyzing the images produced by the sensor and extracting features that can be used for process control. Such customized image processing and analysis capabilities have been developed by Automatix on standard computer platform based machine vision systems. These developments are being pursued under two SBIR programs sponsored by NASA and DoE.

The DoE sponsored SBIR program focuses on welding process monitoring and control. Under this effort, Control Vision is developing a next-generation vision sensing system that is more compact and can be readily interfaced to a

robot or other automated welding equipment, whereas Automatrix is developing vision processing techniques for image enhancement and image analysis for the detection of important weld features, such as the weld seam, molten metal puddle, or keyhole and the calculation of relevant dimensional measurements, such as seam-to-puddle offsets or puddle geometry.

The NASA sponsored SBIR program focuses on thermal spraying process monitoring and analysis. Special techniques have been developed to suppress the intense light of the flame or plasma and to allow the visualization of the powder particle flow carried by the flame or plasma with particular emphasis placed on low-pressure or vacuum spraying processes inside chambers. In conjunction to these sensing developments pursued by Control Vision, Automatrix is developing image processing and analysis schemes for the automated extraction of quantitative process measures such as particle distribution, velocity, and flow rates from the video images produced by the sensor. Figure 1 schematically depicts the viewing system used in conjunction with the image processing and analysis system for computer-assisted visual monitoring of thermal spraying.

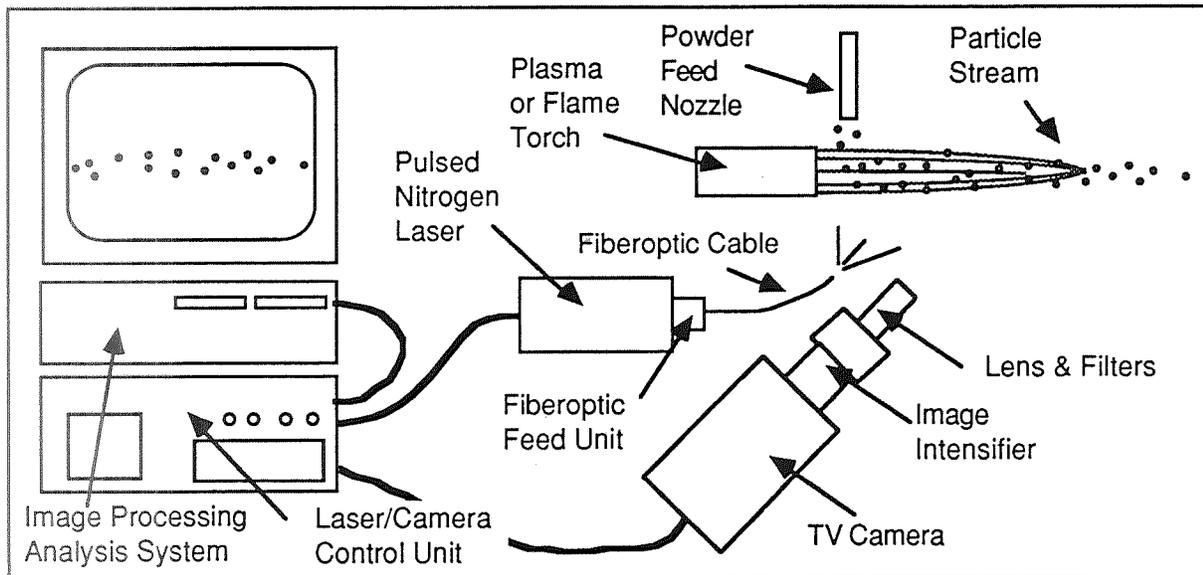


Figure 1 Typical experimental setup for thermal spraying process monitoring

In both welding and thermal spraying, the developed vision sensing and processing schemes can be used both during process development for parameter selection and process understanding or modeling as well as during production for real time process monitoring, process alarming, and ultimately process control (Figure 2).

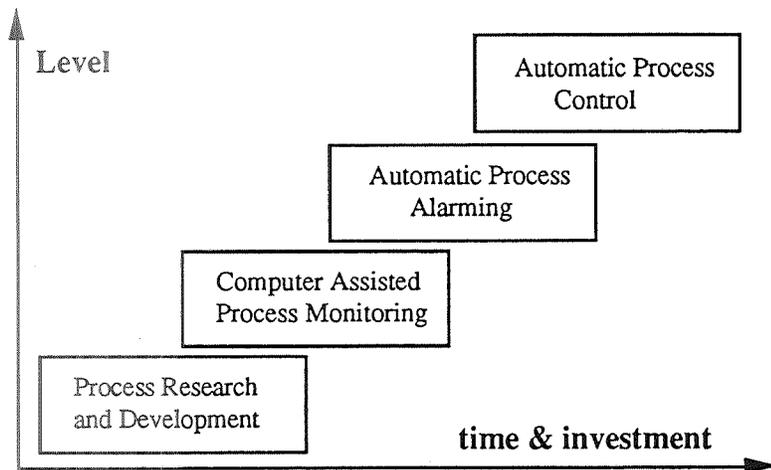


Figure 2 Different levels of application for the proposed monitoring systems

Specifically, such systems can be used during process research and development to evaluate the effect of different process parameters, explore alternative hardware designs, allow for better record keeping during process experimentation, and validate analytical process models. During actual production, these systems can be used for computer-assisted monitoring by computing different process measures and providing this information to an operator through graphical displays, for automatic alarming by comparing these computed measures to expected values and tolerances, and ultimately for real-time process control where the process parameters are automatically adjusted on the basis of the process feedback.

VISION SENSING

As mentioned above, the patented Control Vision viewing system incorporates external illumination in the form of intense pulsed laser light to overcome the extreme variation in scene brightness created by the flame or arc. The laser energy is transported to the welding site through a fiberoptic cable. A xenon flash lamp has also been used as a source of intense pulsed light. The system is also equipped with a narrow-band optical filter to match the laser wavelength and further suppress the arc lighting. The net combination of both temporal and spectral filtering results in a video image that is free of all of the adverse arc lighting effects.

Welding Applications

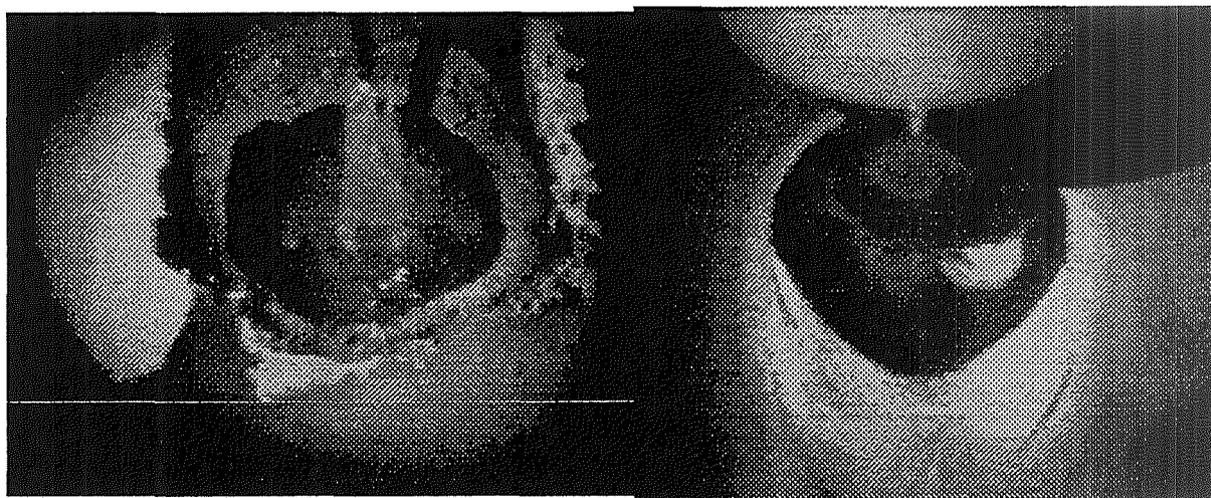


Figure 3. Gas Metal Arc welding viewed using the stroboscopic vision sensing approach

Figure 4 Gas Tungsten Arc welding viewed using the stroboscopic vision sensing approach

The viewing system has been applied to a variety of conventional welding processes including Gas Tungsten Arc Welding (GTAW), Plasma Arc Welding (PAW), Variable Polarity Arc Welding (VPPAW), Gas Metal Arc Welding (GMAW) as well as high energy welding processes such as Electron Beam Welding (EBW) and Laser Welding (LW). The VPPAW and GTAW processes are of particular interest to NASA since they are extensively used for the production of flight hardware such as the Space Shuttle External Tank and Main Engine.

Phenomena not easily seen through a conventional viewing system are readily visualized using the stroboscopic viewing approach making it a unique tool for welding process R&D. Figure 3 clearly shows metal droplet transfer, weld puddle depression, spatter formation, cathodic cleaning, and related phenomena observed during GMA welding of an Aluminum-Bronze alloy. The three frames of Figure 5 provide a clear back side view of the keyhole drilled by the plasma during VPPA welding. No evidence of the intense plasma remains in the image. The viewing system can be used for process monitoring and help a supervisory operator detect setup problems as in Figure 5(b) where the keyhole has moved away from the weld seam or process stability problems such as the unstable keyhole formation of Figure 5(c).

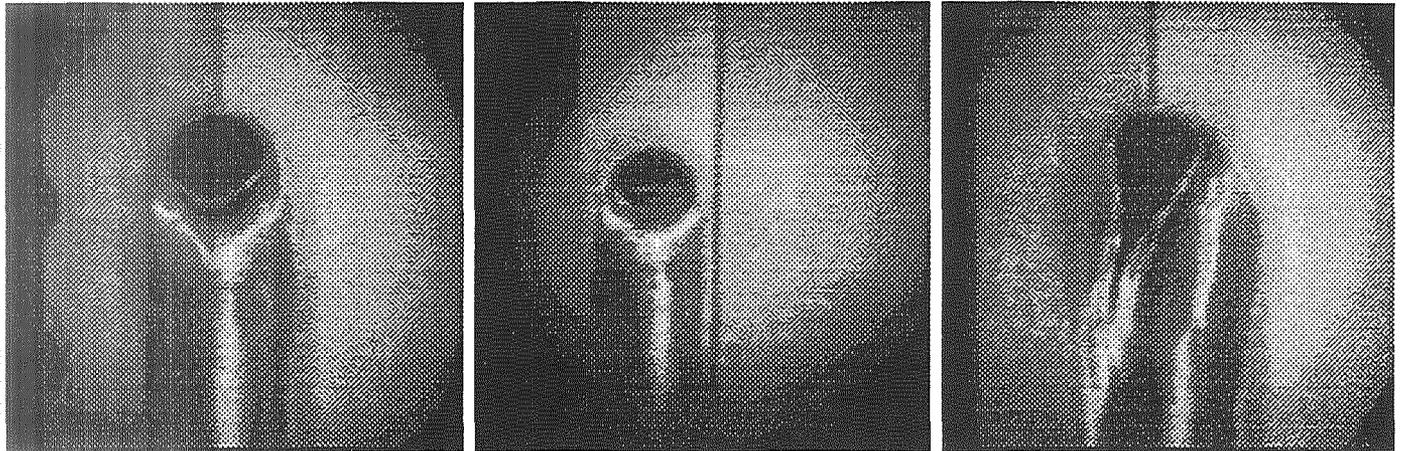


Figure 5 Backside view of the keyhole drilled by the plasma during VPPA welding

Thermal Spraying Applications

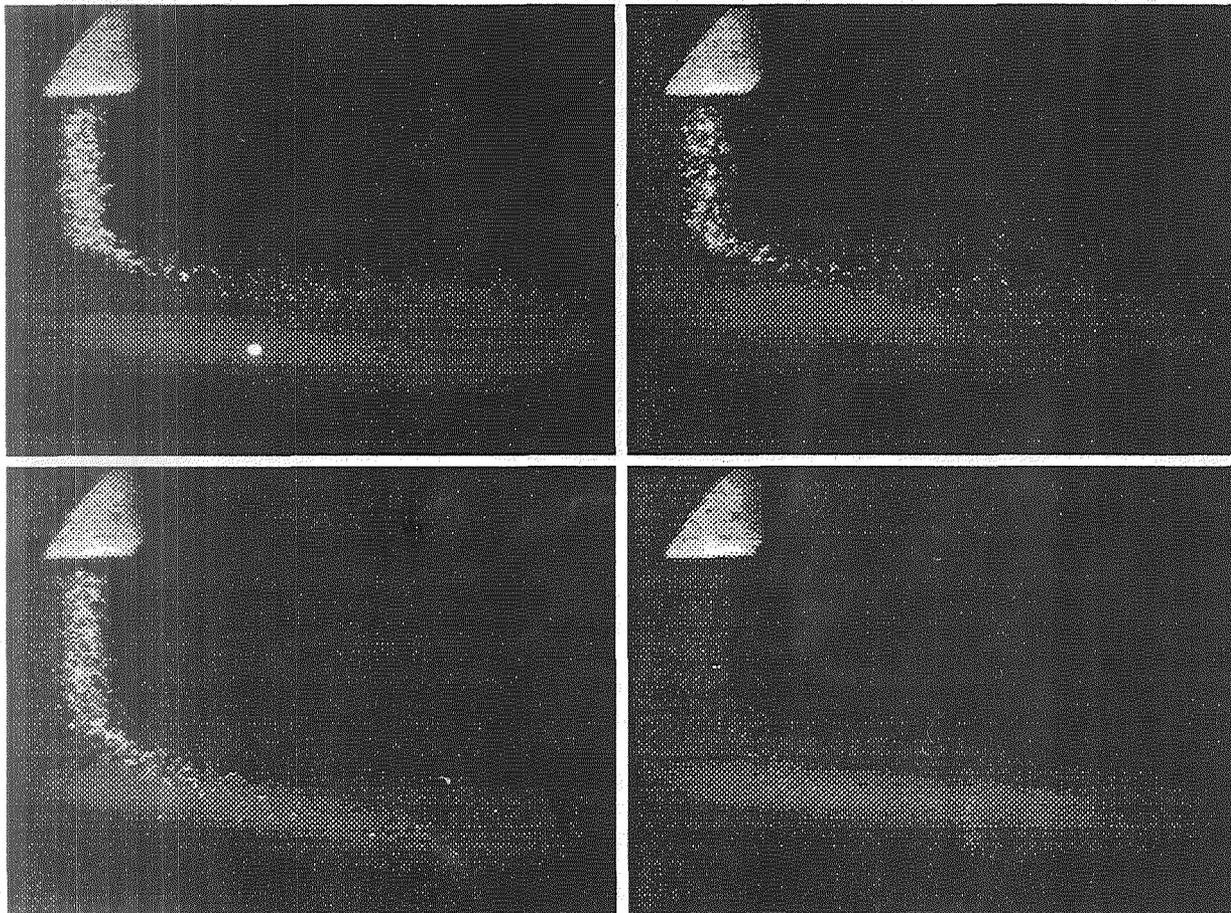


Figure 6. DC plasma spraying showing the effect of varying powder feed parameters in the particle distribution within the plasma plume. As the process parameters are changed the particles ride on the plasma plume (top left and right), overshoot (bottom left), and are optimally carried by the plasma (bottom right).

The vision sensing technology presented in this paper has been used in other manufacturing processes involving a plasma, arc or flame. A family of such processes of particular interest to NASA and the aerospace industry is

thermal spraying which is widely used today for spray coatings and spray forming applications. Thermal coating processes (flame powder, flame wire, wire arc, plasma, etc) utilizing a heat source (chemical combustion or electric arc) to generate an intense flame, arc, or plasma which in turn is used to heat to a molten or elevated temperature solid state, accelerate them to a high speed, and carry and deposit them on a workpiece. Such coatings can minimize or impede environmental effects, improve wear resistance, develop abradable seals or thermal barriers, allow worn-part build-up, or control tolerances. Thermal spraying is also used in forming applications whereby molten material is sprayed layered and rapidly solidified onto a semi-finished product. This approach reduces the number of steps required to arrive at a finished product, compared to more traditional processes. The resultant material also exhibits superior micro-structures by eliminating macro-segregation and promoting finer micro-structures.



Figure 7. Active viewing of twin-wire arc spray in the immediate area of the gun nozzle

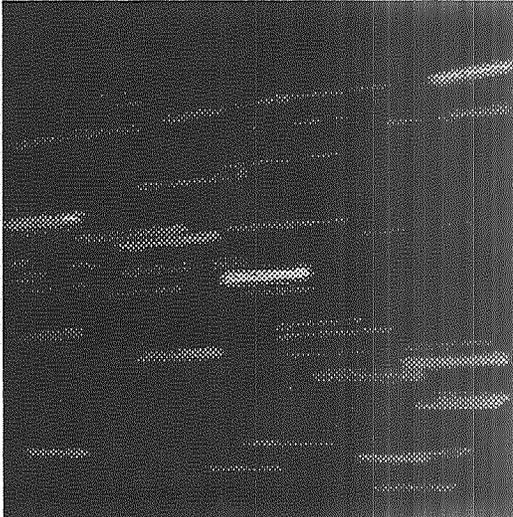


Figure 8 Passive viewing of wire arc spray downstream from the gun nozzle

The strobing parameters can be adjusted so that the plasma can be completely removed from the image. This is demonstrated in Figure 7 showing twin wire arc spray process in the immediate region of the gun. At the opposite end, the laser strobing can also be removed for completely passive viewing. This is only meaningful in the region away from the gun where no evidence of the flame exists. Passive viewing of the plasma is shown in Figure 8. Note that the particles now are shown as light streaks corresponding to the short path they traverse during the sensor integration time. The length of the streaks is thus proportional to the velocity of the particles.

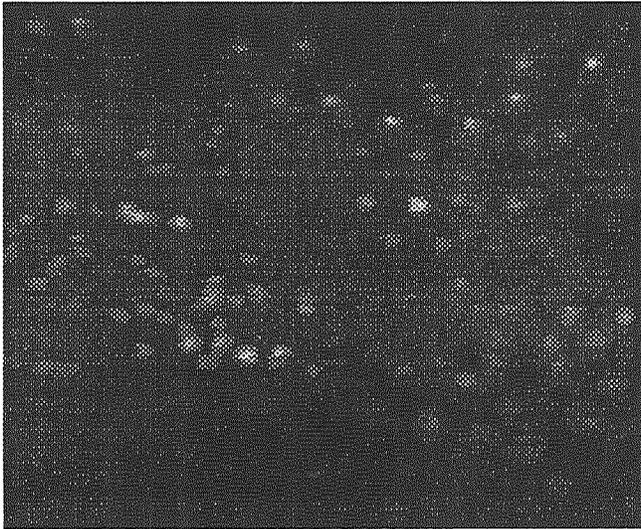


Figure 9 Twin images obtained using two lasers fired sequentially.

Another approach for visualizing the velocity field that has been explored is the use of two lasers or strobes fired in rapid succession. This results in two particle images shifted with respect to one another by an amount that is proportional to the local direction and size of the velocity vector. This is depicted in Figure 9.

Other thermal spraying processes that can benefit from the unique capabilities of the sensor include single wire arc spraying, PTA hardfacing, and gas atomization. Figure 10 shows two digitized images from a videotaped single wire arc sequence.

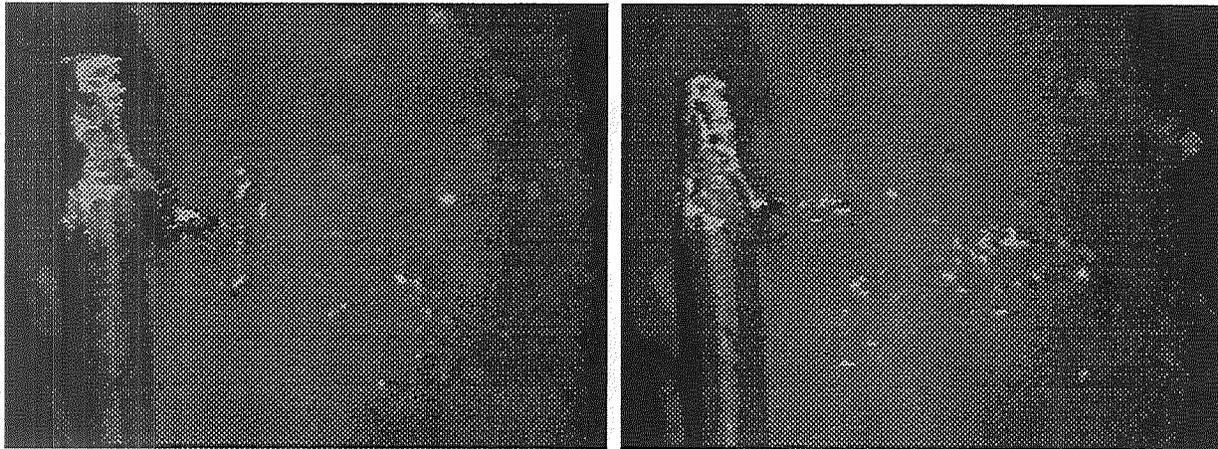


Figure 10 Single wire spraying images

VISION PROCESSING

Under the previously mentioned research programs, Automatrix is developing image processing and analysis approaches for processing and analyzing the video images produced by the sensor. These capabilities are developed on commercial machine vision systems implemented on industrially hardened standard computer platforms using the vision application development environment highlighted in [Schurr 1991].

Image Processing for Image Enhancement

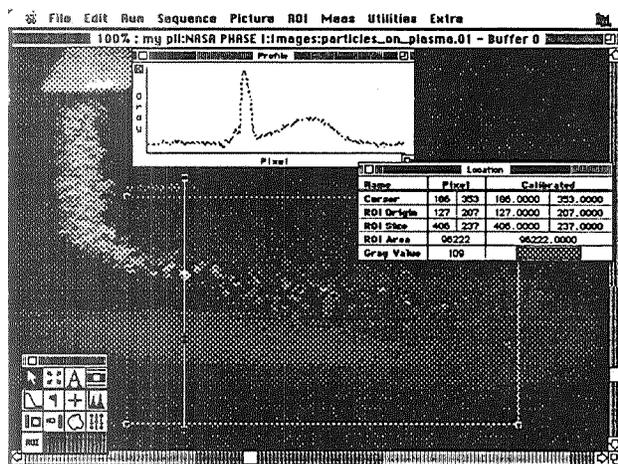


Figure 11 Unprocessed image showing remaining evidence of the plasma plume

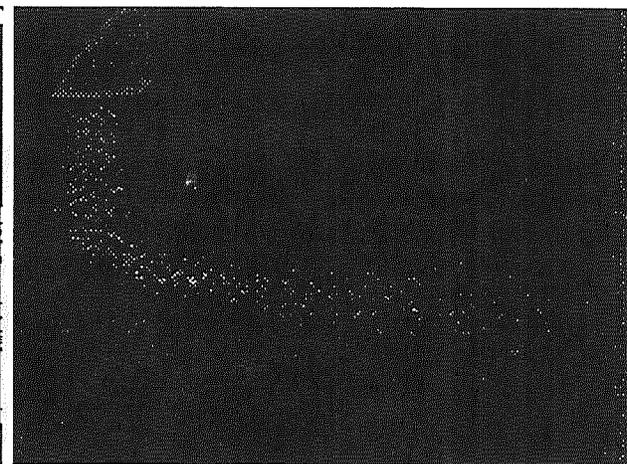


Figure 12 Sample image processing operation to remove the plasma plume

Image processing refers to operations that produce a new image from an original image so as to enhance the subjective image quality, reduce noise, accentuate features of interest, or eliminate image formation problems. Such processing may be useful to make it easier for a human operator to interpret the image or as a preprocessing step prior to any subsequent image analysis. Image enhancement or restoration may also be necessary when analyzing a

recorded sequence of weld images. In such a case, it is no longer possible to change the viewing or picture taking parameters.

For example in thermal spraying applications, the main objective of image enhancement is to process captured images and clean up any remaining evidence of the plasma and/or isolate the particles or the plasma plume. This separation of the particles from the plasma can be based on the fact that the plume has a lower spatial frequency content than that of the particles. This difference is clearly demonstrated by looking at the intensity distribution through the particles and the plasma (Figure 11). Figures 12 and 13 show examples of successful image processing to isolate the powder particles.

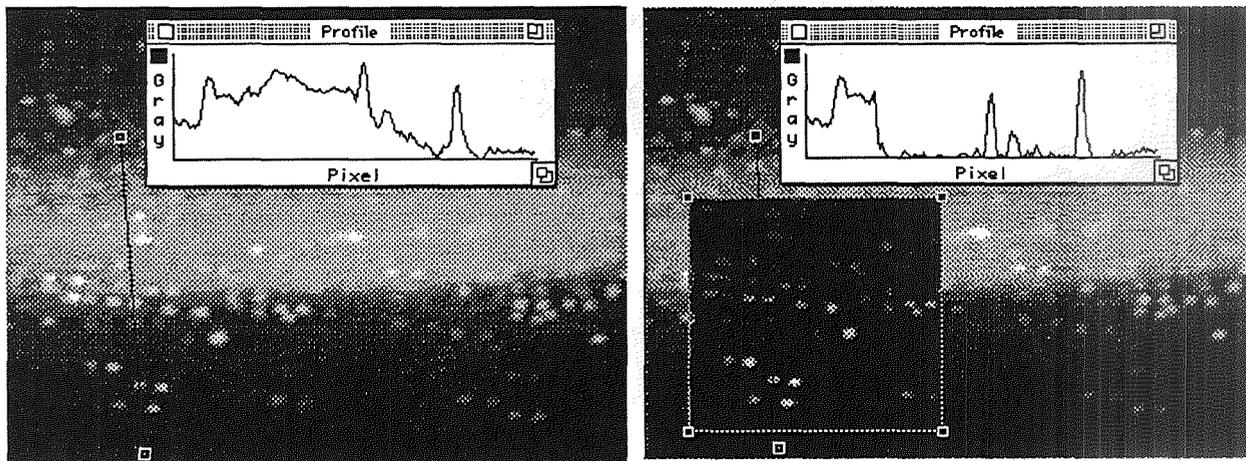


Figure 13 Another example of image processing to remove the plasma plume from the image.

Image Analysis

The main objective of image analysis is to extract features of interest and other quantitative measurements from an image. In the context of welding applications, image analysis is used to:

- to detect the weld joint ahead of the welding torch (its location and size can be used for real-time torch guidance, pre-weld joint inspection, and real-time welding process parameter adjustment) [Agapakis 1990];
- to detect the molten metal puddle or keyhole under the arc (its shape and size can be used for process monitoring and real-time process control);
- to detect the solidified weld bead behind the arc and puddle (its shape and dimensions may be used for post-weld inspection or process control);
- to detect other features of potential interest such as the welding torch, electrode, filler metal wire, molten metal droplets, or the welding arc and plasma.

In thermal spraying applications, image analysis is used:

- to detect the particle detection through image segmentation
- to compute particle area and other particle characteristics
- to compute measures of particle spatial distribution
- to determine the mean particle path
- to compute the particle velocity distribution
- to extract the geometry of the plasma plume or flame

Examples of results of such image analysis for welding applications are shown in Figure 14, where both the centroid of the detected keyhole or puddle and the centerline of the detected joint seam are shown graphically. The distance between the puddle/keyhole centroid and of the seam centerline can be computed during the process and used for seam tracking.

Sample results of image analysis for thermal spraying applications are shown in Figures 15 and 16.

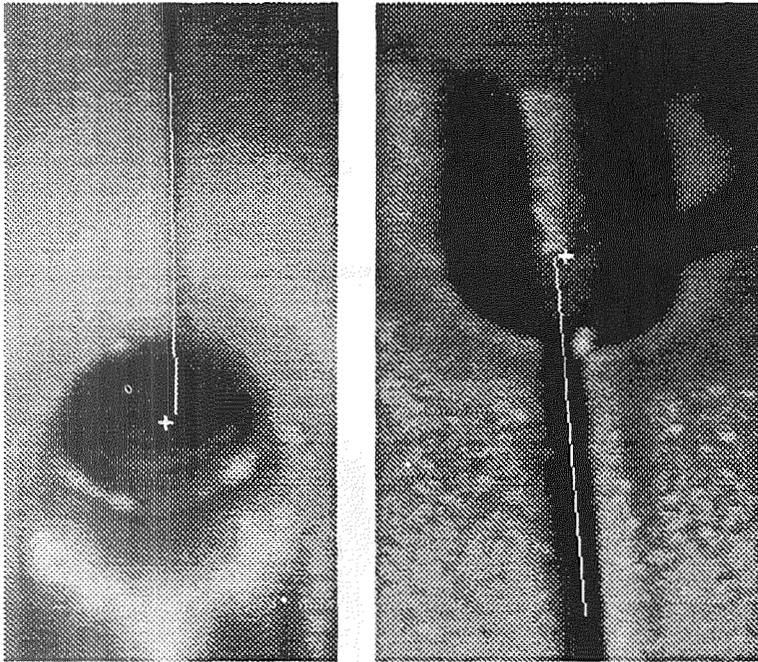


Figure 14 Results of image analysis during VPPA and GTA welding demonstrating the detection of the seam centerline and of the keyhole or puddle centroid.

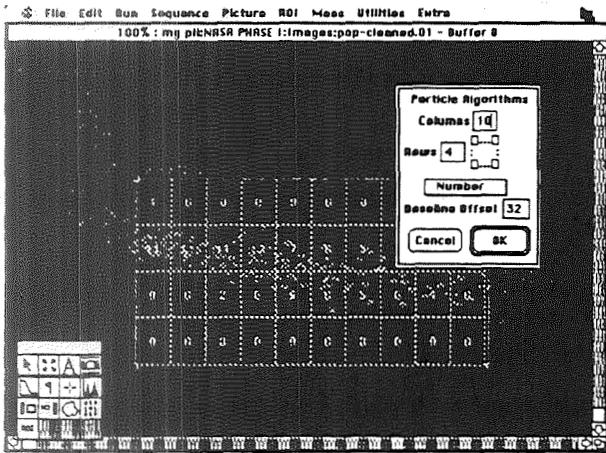


Figure 15 Determination of particle spatial distribution

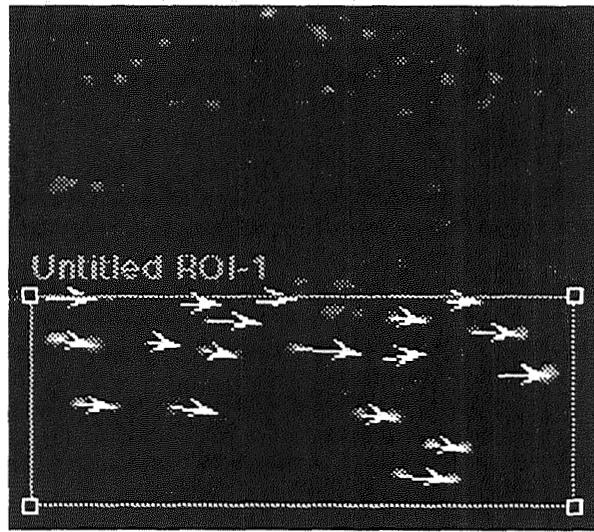


Figure 16 Determination of local velocity vector field

ACKNOWLEDGMENTS

The work presented in this paper has been partially sponsored by a Department of Energy Phase II SBIR grant and a NASA Phase I SBIR contract. In addition to the authors, Richard Lu, Phil Lamoreaux, and Thad Hoffman, and John Lagerquist have contributed to these research efforts.

REFERENCES

- Agapakis, J.E., Katz, J.M., Friedman, J.M., and Epstein G.N., "Vision-Aided Robotic Welding: An Approach and a Flexible Implementation" *International Journal of Robotics Research*, Vol 9, No 5, October 1990, pp 17-34.
- Schurr, G., "A Multi-level Application Development Environment for Machine Vision System Integration, proc. *International Robots and Vision Conference*, Detroit, MI, October, 1991.

ROBOTIC VARIABLE POLARITY PLASMA ARC (VPPA) WELDING

Waris S. Jaffery
Boeing Defense & Space Group
Huntsville, AL 35824

N 93-25675

150514
P-10

© 1992 The Boeing Co. All Rights Reserved

ABSTRACT

The need for automated plasma welding was identified in the early stages of the Space Station Freedom Program (SSFP) because it requires approximately 1.3 miles of welding for assembly. As a result of the Variable Polarity Plasma Arc Welding (VPPAW) process's ability to make virtually defect-free welds in aluminum, it was chosen to fulfill the welding needs. Space Station Freedom will be constructed of 2219 aluminum utilizing the computer controlled VPPAW process. The "Node Radial Docking Port," with its saddle shaped weld path, has a constantly changing surface angle over 360° of the 282 inch weld. The automated robotic VPPAW process requires eight-axes of motion (six-axes of robot and two-axes of positioner movement). The robot control system is programmed to maintain Torch Center Point (TCP) orientation perpendicular to the part while the part positioner is tilted and rotated to maintain the vertical up orientation as required by the VPPAW process. The combined speed of the robot and the positioner are integrated to maintain a constant speed between the part and the torch. A laser-based vision sensor system has also been integrated to track the seam and map the surface of the profile during welding.

INTRODUCTION

The Space Station Freedom Program (SSFP) is an American-led international project. National Aeronautics and Space Administration (NASA), with subcontractors, will develop, launch and operate a permanently-crewed base in space by the mid-1990's. The space station will consist of an external structure, larger than a football field, connected to pressurized modules in which men and women will live and work in space. Boeing's role is to build the habitation and laboratory modules.

In order to manufacture the SSFP modules, new processes have to be developed. Innovative approaches have to be taken which will provide certifiable processes with the highest quality of workmanship in every step of the production of Space Station Freedom.

After careful review of all the process requirements, design configuration of parts and with flexibility in mind, the decision was made to acquire a welding system which will allow maximum flexibility and programmability. The welding system is to accommodate changes in configuration as the weldment designs are refined, also the workcell can be easily adapted to weld assemblies other than the Node Radial Docking Port (Figure 1).

SYSTEM SPECIFICATION

A specification was written to establish the basic requirement of the system, which will not only fulfill the process requirements, but also allow room for future expansion. In order to determine the robot's arm reach, motion simulation software was used.

Utilizing Deneb Robotics, Inc.'s graphic simulation package IGRIP (TM), the VPPA weld process profile motion was simulated. The simulation provided the preliminary location of the robot and weld fixture, which led to determination of the robot's arm reach. After performing several simulations, it was determined that a system with a minimum of 100 in/250 cm reach would be needed to perform the complex weld task. The simulation work was done in conjunction with NASA's Metal Processes Branch of Marshall Space Flight Center (MSFC) Productivity Enhancement Center.

Specifications were sent to nine different vendors. A vendor briefing was held at Boeing to answer all the questions regarding the specification. Five vendors participated in the question and answer session. Later, three vendors responded with proposals.

PROCESS

B. P. VanCleave, of The Boeing Company, along with Hobart Brothers of Troy, Ohio, combined the Plasma Arc process with Variable Polarity features in the late 1960's. As a result of VanCleave's promising work, a VPPAW research and development project was initiated in the late 1970's at NASA's Marshall Space Flight Center (MSFC) located in Huntsville Alabama, to determine the potential usage of the VPPAW process in the fabrication of space programs. Based on the success of the VPPAW process on the Space Shuttle's External Tank, Boeing selected the process as a key part of manufacturing the Space Station Freedom's Module and Node structures. The uniqueness of the VPPAW process is in the controlled manipulation of weld current pulse width and the addition of a controlled pulse of reverse polarity that supplies the cleaning action. The high velocity Argon gas flowing through the center of the orifice is an essential element to the process. The high velocity of the plasma arc punches a hole through the work piece known as a "Keyhole". The creation of the "Keyhole" in the workpiece opens up ahead of the plasma jet in the upward direction (direction of travel), and an impurity free molten puddle of aluminum flows downward, and solidifies behind the plasma jet. Because of this "blow-thru" mode of welding, the torch must maintain a "vertical-up" position during the entire welding process. Maintaining this orientation for joining straight or cylindrical plates is relatively simple and has been achieved using simple tools such as a two-axis manipulator and stationary fixture. Maintaining the "vertical-up" orientation becomes a complicated task when the weld seam geometry is in a three-dimensional orientation.

THE SYSTEM

The automated robotic VPPAW system, designed and integrated per Boeing's specifications by Cincinnati Milacron (currently ABB Robotics, Inc.) of Greenwood, South Carolina. The schematic diagram of the VPPA welding system is shown in Figure 2.

System Description

The cell is comprised of a Cincinnati Milacron T3(R)-696 Robot with control system, a Hobart HAWCS(R) VPPA Welding system, two Aronson RAB-60(R) positioners, a Stenning Weld viewing system, a perimeter guarding system, a gas manifold and monitoring system, and an operator console. The perimeter guarding system divides the workcell into three zones. In zones A and C, positioners are located, while zone B is the robot HOME zone (Figure 3). The two positioners were opted for the system to increase the efficiency of the workcell. While one assembly is being welded in one zone, a second assembly can be prepared for welding in a second zone.

Control

Cincinnati Milacron Acramatic(R) Version 5i robot control utilizes distributed microprocessor architecture and an MS-DOS(R) base operating system. It consists of a large option of hardware/software capabilities designed to facilitate efficient interfacing and controlling of external equipment and processes. In this workcell, the controller is interfaced with the safety monitoring system and weld control system. The robot control system acts as the master controller. A teach pendant is used to teach the robot and for programming the controls. The pendant could also be used for making correction of taught points for part deviations from the engineering design. Off-line programming capability also exists.

Robot

The robot is a standard six-axes articulated all electric servo-controlled unit consisting of dual resolver feedback for absolute positioning capability. The robot has a 275 lbs. load capacity at the end of the arm. End of the arm tooling consists of a welding torch with cross slide, AVC, torch rotator, wire guide manipulator and viewing camera (Figure 4). The position accuracy is within 0.010 inch.

HAWCS(R) VPPA Welding System

The HAWCS(R) system consists of an 80386 AT(R) industrial computer. The HAWCS(R) is the heart of the welding system. All the welding parameters and process controls are programmable. The typical welding program sequence is shown in Figure 5. It is the sequence of instructions that is executed by the controller to control the entire weld procedure as programmed. The automated voltage control (AVC) provides feedback to the control system, through utilization of computer software, and arc voltage is adjusted by changing the torch-to-work piece

distance accordingly. All the procedure editing and execution, as well as system configuration and administration, are accessible through a VGA monitor which has a touchscreen interface for data entry. Each weld procedure is comprised of four segments: Start Up, Main, Termination and Emergency. The teach pendant provides remote display of weld parameters as well as trim capability, gas purge control and weld start/stop control. The HAWCS(R) control communicates with the robot control via its digital input/output (I/O) control to Robot I/O control. Position and velocity data is communicated to the HAWCS(R) control by a programmable digital to analog convertor which emulates a single channel encoder by sending a five volt DC digital pulse every 0.010 inches of torch center point motion. The welding parameters can be printed as well as stored on the disk during the weld in real time.

Operator Console

The operator console contains the operator interface for selecting programmed weld sequences, perimeter guarding, system status and intervention. Also included at the console are the weld viewing system, pendant for the control of the robot and for the HAWCS(R), and the monitors for seam tracking and bead profiler system.

Positioner

The RAB-60(R) is a two-axis positioner with servo-motor which is controlled by the robot control. The positioner provides tilt and rotation axes for part positioning. The positioner can carry a payload of up to 6,000 pounds. The positioners can rotate $\pm 458^\circ$ and tilt 0 to 110° . The positioner axes may be moved using the robot pendant in either Manual Mode or in Teach mode. The positioner axes can also be programmed for part set-up in one zone while the robot is operating in Auto Mode in another zone.

Weld viewing system

In order to provide a safe robotic weld workcell and also to detect the shape of the groove and view the torch and wire tip manipulator, a remote viewing method has been utilized. A viewing system with filtration, which provides a uniform image of the welding arc, the weld pool, and the adjacent area was selected. The remote controlled viewing system provides substantially more real-time information to the welding operator and aids in making the needed adjustment while the weld is in progress.

The Stenning viewing system is comprised of a camera and light source located at the torch and a monitoring and control unit at the operator console. The viewing system has remote controlled focus, filter and iris, it also has the capability to video tape the weld utilizing a VCR.

Gas control

The gas manifold system is for dispensing and metering shielding and plasma gas. Manifold pressure is monitored by a pressure switch which is interfaced with safety control system which runs as a shell to the entire system. When pressure drops below 600 psi the LOW PRESSURE light is illuminated on the operator console.

SEAM TRACKING AND BEAD PROFILING

Seam tracking and bead profiling sensors have been added to the robotic end effector. It is intended to provide fully automated VPPAW welding with override capability of the remote control functions of cross slide and torch rotation control at the operator console. A square butt weld joint configuration with a slight edge chamfer will be utilized for the first pass seam tracking. The chamfer will not be visible for second pass welds, which also must be tracked. The system will be used to track the seam during root pass welding, as well as cover pass welding. For second pass tracking, the system will use the encoder tracking information from the first pass combined with a template matching algorithm based on the data from the first pass. The system mounting is shown in Figure 6.

Seam Tracking

A non-contact sensor will align the torch with the seam to be joined. It will position and track the center line of both the root and cover passes. The sensor is mounted on a position relative to the torch so the seam area is just ahead of the torch center point. The sensor moves forward along the seam in tandem with the weld torch,

scanning the seam area from the torch side of the weld joint, and mapping the surface profile so the complete and true geometry of the seam is known and the torch is able to anticipate it and adapt accordingly. The sensor has been mounted in such a way that the seam is always in the field of view.

Bead Profiling

A non-contact laser vision sensor mounted behind the torch will be used to measure all of the weld parameters as shown in Figure 7 (at end of report). The prime purpose of the bead profiler is to measure and control the solidified weld bead parameters during the VPPA weld processes. The system will operate during VPPA welding and will be used to correct asymmetric bead profile. The sensor has been mounted one inch below the torch and measures the weld bead as soon as possible after solidification. The bead profiler will also control rotation of the torch to correct the bead's geometry by eliminating asymmetrical undercut contour. As the bead profile sensor senses the bead undercut, or any other asymmetrical contour of the bead, the torch would be rotated to correct this.

The data collected from seam tracking and the bead profiler's geometry parameter acquisition will be used for real time adaptive control and/or statistical process control (SPC) analysis.

Safety

The Perimeter Guarding System is comprised of six photo cell beam stands located in the front and rear of the cell. Each beam stand is comprised of three photo cells (sender / receiver). The photo cell beams guard against inadvertent entry into the cell or an unauthorized area of the cell. There are two gates located in zone "A" and "C" Figure 3. The gates are interlocked to the perimeter guarding system by use of a unique configured key system. The operator must have permission to enter that area of the cell before the key can be removed from its holder in the operator console. The operator can override the guarding system, under certain given conditions, allowing access to the non-working zone of the robot. Any unauthorized entry into the working zone of the cell results in an Emergency Stop being generated.

SYSTEM OPERATION

In order to run the system, weld parameters and robot motion have to be programmed. Several standard sub routines have been developed in the system for weld parameters and robot motion, thus minimizing the programming efforts required for different parts.

The automated welding system was installed in January, 1991. After initial training, the system became operational in February, 1991. The weld schedules for different thickness plates were developed and tested utilizing vertical test stands. Off-line simulation was used to determine the optimum location of the positioner in zone "A" and "B" Figure 3, to place the saddle shaped node radial docking port, window panel and cupola within the robot's reach. The node radial docking port was the first part to be welded on this system.

Part programming

The VPPAW process requires that the torch remain basically perpendicular with a 3° lead angle with respect to the part to maintain optimum molten metal orientation. The robot must be programmed to maintain Torch Center Point (TCP) orientation perpendicular to the part while the part is tilted and rotated to maintain the vertical-up orientation. The combined speed of the part and the torch must remain constant during the weld, with the exception of the start and termination points. The start-up requires ramp-up acceleration, while the termination requires tapered deceleration.

The node radial docking port weld path is saddle shaped, and thus has a constantly changing surface angle over 360°. The weld fixture and node radial docking port must be placed on the positioner with a clamp ring bolted to the fixture. The clamp ring is very carefully placed to assure that no gap exists between the part and the fixture. The clamp ring must be tightened down following a star pattern tightening sequence to ensure that there is minimum distortion of the part. The clamp ring also provides a heat sink function.

The first programmed point was generated by tilting and rotating the positioner and moving the robot to the weld starting point so that the torch and part surface became perpendicular to each other. Subsequent points were generated by tilting and rotating the positioner about 3° while moving the robot vertical up so that the part remained

perpendicular to the torch. Approximately 130 points were generated along the saddle path. In order to program the saddle path accurately, devices such as a trisquare, inclinometer and a spring loaded torch head were utilized. At each point, the trisquare was used to ensure that the part was vertical with respect to the earth and the torch was perpendicular to the part on the horizontal axis. The inclinometer was used to ensure that the torch orientation was such that the tip of the torch was 3° above the horizontal plane. The spring loaded torch head was used to make sure that the torch tip would be on the center of the seam. Each point was measured and checked carefully prior to programming.

The total weld path motion was 282 inches, requiring 200 + points and eight-axes of coordinated motion were programmed to accomplish the saddle weld. Three different welds were made on the part, tack at eight inches per minute (IPM), root pass at seven IPM and cover at ten IPM. The welds were right on the seam during the entire weld procedure. The robot tracked the seam very accurately with minimal (+/- 0.010") deviation from the weld path. One concern was that there would be part deflection due to heat from the weld which would effect the path; however, there was very minor deflection of the part, due to the heat sink of the clamping ring and fixturing.

CONCLUSION

The production workcell has been installed and successful VPPA welding has been accomplished on development articles. The cell is currently operational for the production of Space Station Freedom.

The use of the eight-axes controlled motion is the first known production application of a robot in tandem with the Variable Polarity Plasma Arc Welding System. Implementation of this system will allow expansion of the high quality VPPA weld process to part configurations that have not been considered good candidates.

ACKNOWLEDGEMENTS

This work was accomplished under a U. S. Government contract, number NAS8-50000, through the Marshall Space Flight Center, Huntsville, Alabama, for the Space Station "Freedom" Program.

The author wishes to express his gratitude to the MSFC and Boeing personnel for their support and efforts in preparing this paper, and to the salient leadership and directions of all participating Boeing personnel.

REFERENCES

1. Brosilow, Rosalie ed., "Space Shuttle Flies on Computer Welds," *Welding Design and Fabrication*, pp 25-31 (August 1989).
2. Cary, H. and Barhorst, Steve, "Advances in Welding Science and Technology," *Proc of an International Conference on Trends in Welding*, Gatlinburg, Tennessee (May 18-22, 1986).
3. Swinghammer, R. J., "Unique Variable Polarity Plasma Arc Welding for Space Shuttle" (October 1985) NASA TM-86536.
4. VanCleave, B. P., and Gain, W. R., "Keyhole Plasma Arc Welding of Aluminum," *American Welding Society* (October 7-8 1980), pp. 109-122.
5. Jaffery, Waris S., "Automated Robotic Variable Polarity Plasma Arc Welding (VPPAW) for Space Station Freedom Program," *International Conference on Trends in Welding*, Gatlinburg, Tennessee (June 1-5, 1992).
6. Jaffery, Waris S., "Implementing A Robotic VPPA Welding System For Space Station Freedom Program," *SME Conference At IMTS, Robot Welding*, Chicago, Illinois (September 11, 1992).

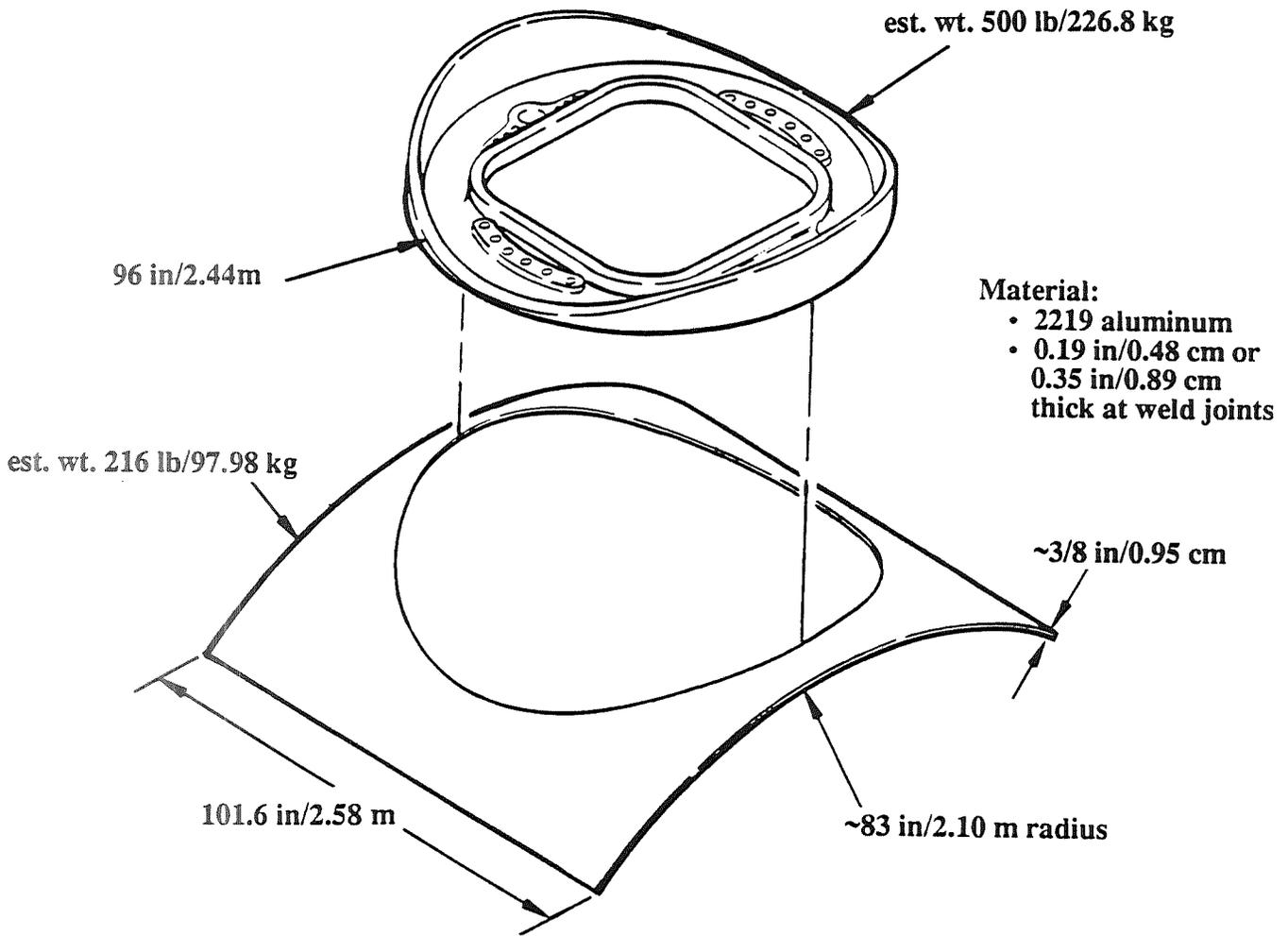


Figure 1. Node Radial Docking Port (Saddle Weld)

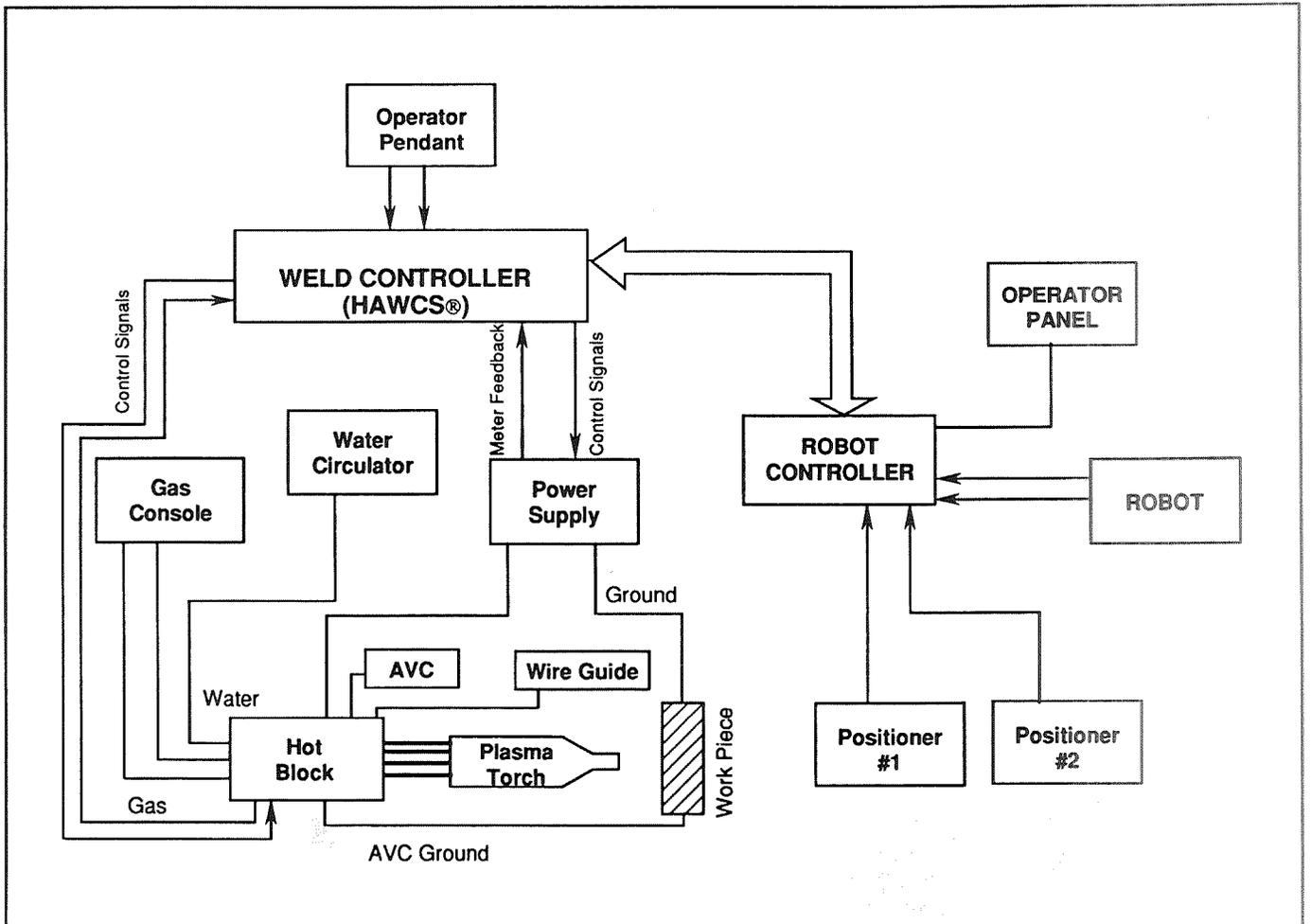


Figure 2. Schematic Diagram of the VPPA Welding System

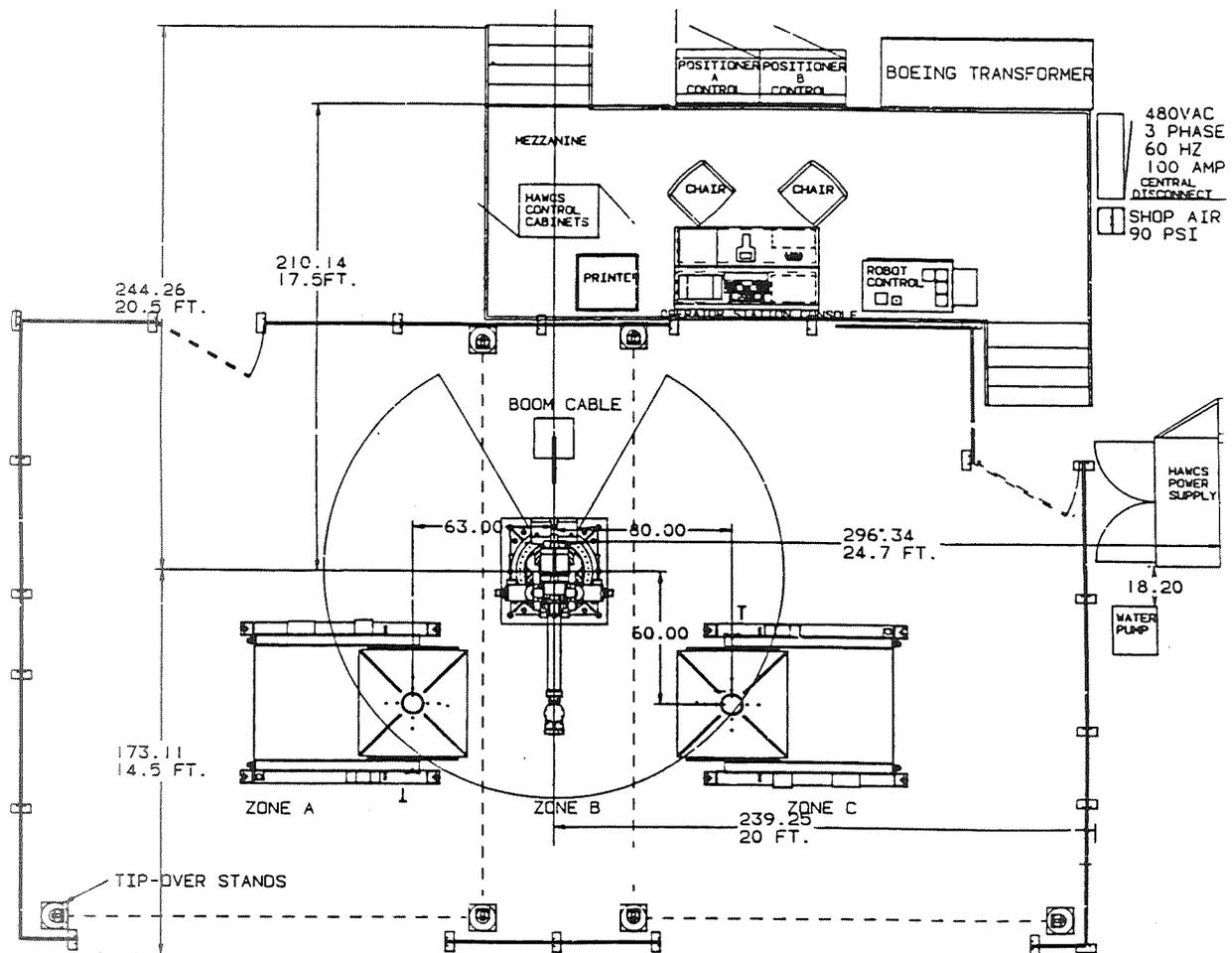


Figure 3. Workcell Layout

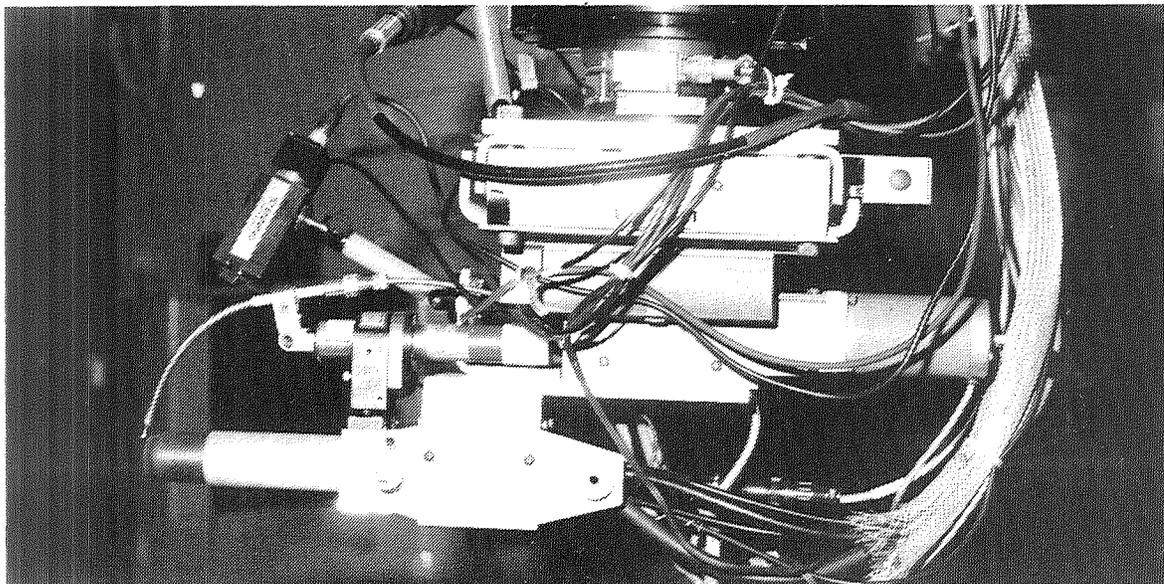


Figure 4. Robot End Effector Tooling

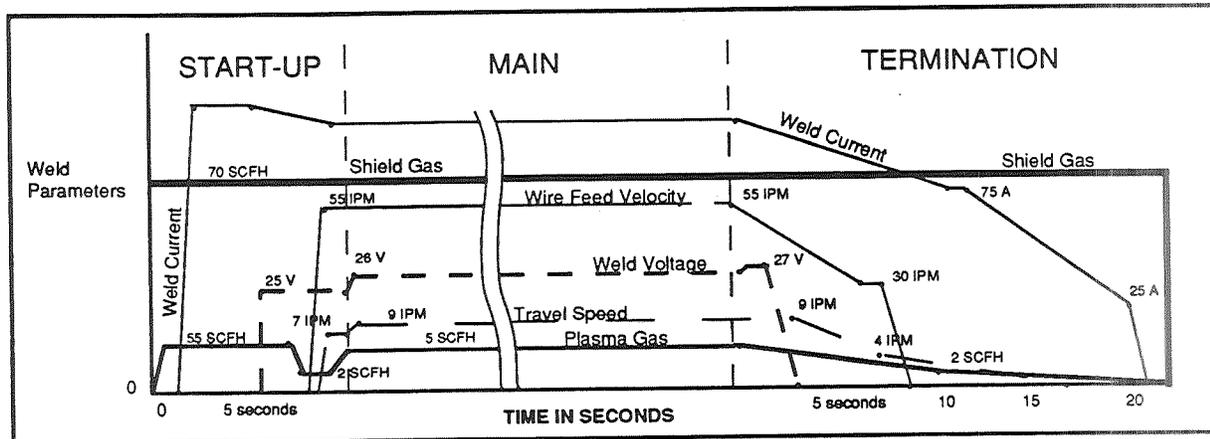


Figure 5. Weld Program Sequence

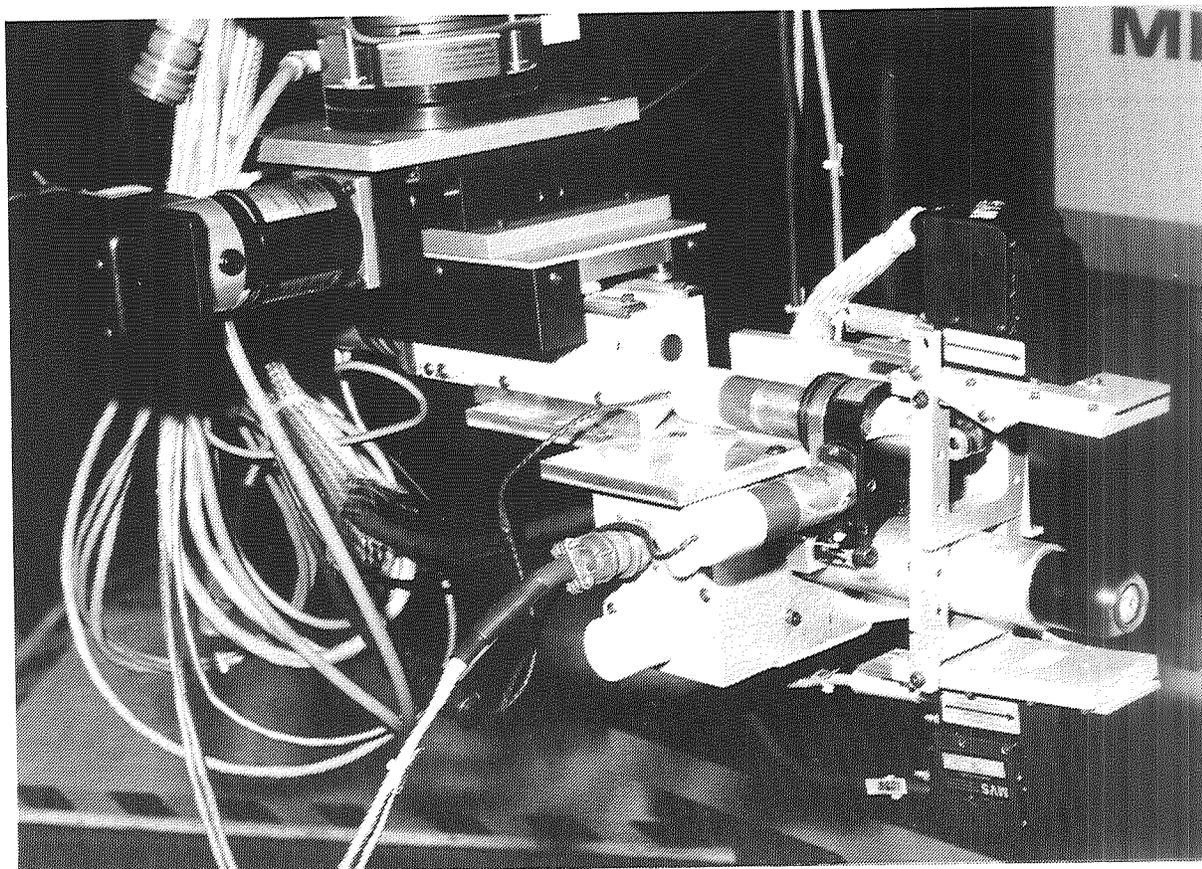


Figure 6. Seam Tracking and Bead Profiling System Mounting

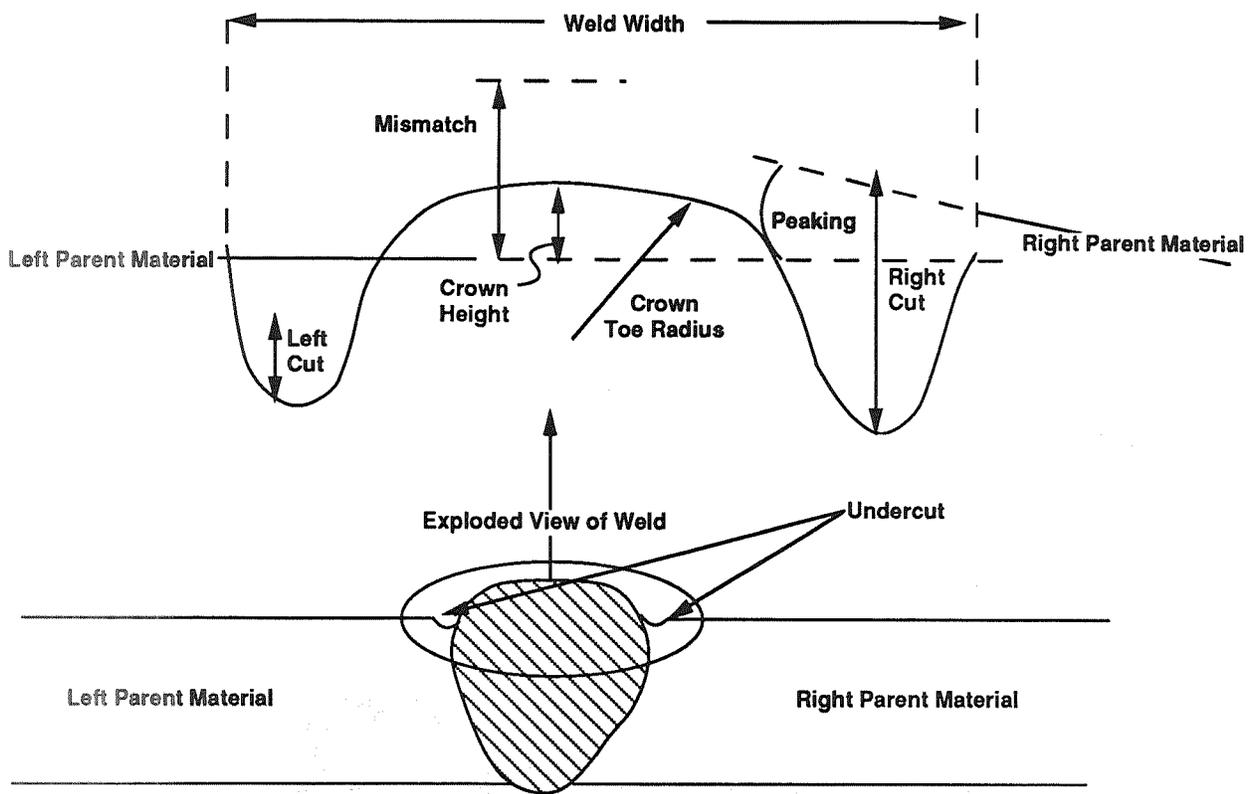


Figure 7. Weld Bead Parameters

FIRMWARE DEVELOPMENT IMPROVES SYSTEM EFFICIENCY

E. James Chern and David W. Butler
Materials Branch / Code 313
NASA Goddard Space Flight Center
Greenbelt, Maryland 20771

345-60

150515

P-10

ABSTRACT

Most manufacturing processes require physical pointwise positioning of the components or tools from one location to another. Typical mechanical systems utilize either stop-and-go or fixed feed-rate procession to accomplish the task. The first approach achieves positional accuracy but prolongs overall time and increases wear on the mechanical system. The second approach sustains the throughput but compromises positional accuracy. A computer firmware approach has been developed to optimize this pointwise mechanism by utilizing programmable interrupt controls to synchronize engineering processes "on the fly". This principle has been implemented in an eddy current imaging system to demonstrate the improvement. Software programs were developed that enable a mechanical controller card to transmit interrupts to a system controller as a trigger signal to initiate an eddy current data acquisition routine. The advantages are: (1) optimized manufacturing processes, (2) increased throughput of the system, (3) improved positional accuracy, and (4) reduced wear and tear on the mechanical system.

INTRODUCTION

Many industrial production processes such as machining, cleaning, assembling, labeling, inspections, etc., require mechanical maneuvering of components or assemblies from one position to another. Typical state-of-the-art mechanical systems use machine tool or robotics motion controllers with stop-and-go or fixed feed-rate mechanisms to meet positioning requirements. The stop-and-go approach achieves positional accuracy but prolongs overall processing time and increases wear on the mechanical system. The fixed-rate sampling approach sustains the throughput but compromises positional accuracy. We have employed a motion controller board and developed software programs to achieve both positional accuracy and sustained throughput. Interrupts generated by the motion controller firmware can be used to synchronize engineering processes "on the fly" for those processes that require minimum or no dwell time at given locations.

The hardware and associated software have been successfully implemented in an eddy current imaging system which consists of a system controller, an eddy current instrument, a mechanical motion controller card, and a two-axis x-y linear table with incremental encoders. The scan routine programs the mechanical controller card with the scan parameters and pre-determined measurement positions. The mechanical controller card's microprocessor constantly compares the real-time probe position from the encoder feedback with the preset acquisition coordinates. An interrupt is generated by the mechanical controller and is used by the system as the trigger signal to initiate the data acquisition routine. Although the concept is demonstrated on an eddy current image system, the interrupt control mechanism can be applied to many other engineering processes.

In this paper, we reviewed hardware and software requirements for implementation of firmware interrupt controls in an eddy current imaging system and other practical applications. Laboratory setup and experimental procedures for the benchmark comparison of the current and conventional approaches using the eddy current imaging system are also described. The results from the parametric experiments clearly demonstrate the improvement in system efficiency. Essential C-language source codes for the interrupt control routines are provided in the Appendix for user reference.

BACKGROUND AND APPROACH

Recent advances in personal computer and electronic technology have enabled the development and operation of various stand-alone concurrent engineering stations. These advancements also facilitated the development of evaluation engineering, nondestructive evaluation (NDE), signal acquisition, image processing, and data presentation techniques. NDE methods are widely utilized in manufacturing and service

industries for quality assurance and related applications[1,2]. Image-oriented data presentation which directly correlates acquired engineering parameters with component coordinates is generally the preferred way of displaying results.

Based on the underlying physical principles of NDE imaging methods, images can be acquired in two forms depending on whether or not the sensor or specimen is manipulated with respect to the other. Imaging systems such as ultrasonic C-scan and eddy current imaging, have to rely on a mechanical scanner to physically maneuver the probe relative to the specimen point by point over the area of interest to acquire images. A typical pointwise NDE imaging system consists of three major components: a system controller to control instruments, command movements, and acquire data; instrumentation to excite the sensor and measure desired signal parameters; and a mechanical scanner to relatively scan the sensor over the area of interest on the specimen. The block diagram and a sketch of a typical eddy current imaging system is shown in Figure 1.

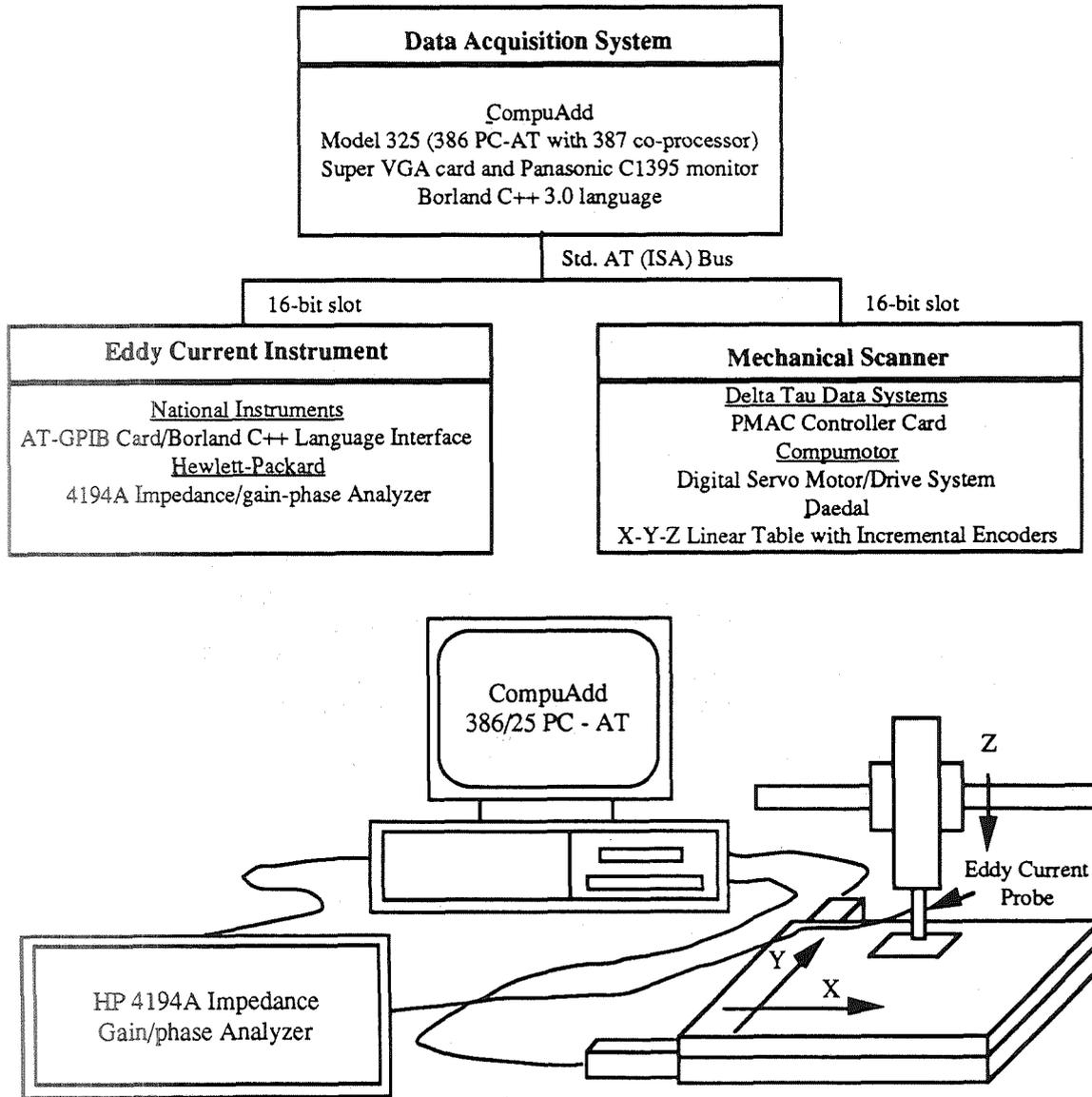


Figure 1. The block diagram and the sketch of the prototype interrupt based eddy current imaging system.

The ideal pointwise imaging system is to command the scanner to scan at a desired speed and fetch measurements at the designed positions "on the fly". However, due to hardware and software constraints, data acquisition is commonly accomplished by either stop-and-go or fixed rate sampling. The principle of the improvement is to utilize the newly available microprocessor based motion controller card as an intelligent controller which initiates and controls the data acquisition process.

The specific approach is to develop firmware routines which enable the motion controller card's microprocessor to constantly compare the real-time probe position from the encoder feedback with the preset acquisition coordinates[3]. An interrupt signal is transmitted to the system controller as a trigger to initiate a data acquisition routine when the positional conditions are met. We have devised a position-driven closed-loop mechanical system for NDE applications. The system uses interrupts generated by the mechanical system at the designed positions, to trigger and initiate the data acquisition routine for the measurements[4].

SYSTEM CONFIGURATION

The improved eddy current imaging system consists a CompuAdd 325 as the system controller, a Hewlett-Packard 4194A Impedance/gain-phase Analyzer as the signal drive and measuring instrument, a Delta Tau Data Systems PMAC motion controller card, Compumotor Plus motors, and a Daedal X-Y linear table with incremental linear encoders as the mechanical scanner. The system controller interfaces with the scanner using the PMAC mechanical controller card through the industry standard architecture (ISA) PC-bus and acquires eddy current impedance data from the HP 4194A through an IEEE-488 interface bus.

The system controller commands the mechanical system such that, the probe traverses the area of interest in a raster pattern. The impedance of the probe is acquired by the impedance analyzer during the scan. Firmware was developed to enable the PMAC to generate interrupts in the system controller as trigger signals to initiate data acquisition sequences. The interrupt structure between host controller and peripheral PIC is shown in Figure 2. The C source code of the interrupt routine is listed in the Appendix.

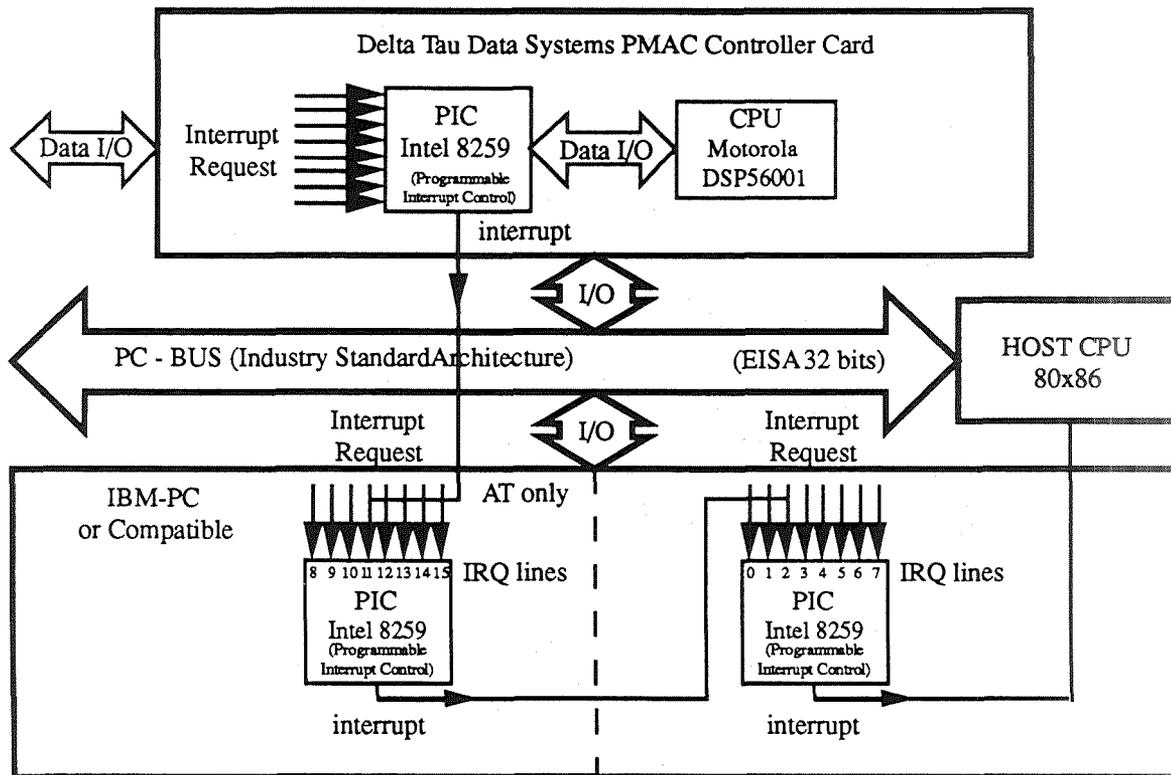


Figure 2. The interrupt structure for the Host PC and peripheral PICs of the prototype interrupt-based eddy current imaging system.

The mechanical scan subroutine programs the PMAC with the necessary scan parameters such as the home position, distances, velocities, and accelerations as well as pre-determined positions where measurements are to be performed. The PMAC microprocessor constantly compares the real-time probe position from the encoder feedback with the preset acquisition coordinates. An interrupt signal is generated by a Programmable Interrupt Controller (PIC) on the PMAC card when the positional conditions are met. This signal is then received by another PIC located in the host controller. This PIC subsequently generates an interrupt in the host CPU. This interrupt is used by the host CPU as the trigger signal to initiate the data acquisition routine and synchronize other events. The PMAC PIC continually generates interrupts until the scan subroutine is completed.

Since the linear encoders are independent of the mechanical drives, the interrupts are generated precisely at the desired coordinates. The only constraint is that the speed of the scanner is limited by the time needed to complete the data acquisition routine and transfer the data through the interface bus. Also, this improvement can only apply to engineering processes that do not require prolonged dwell time. Figure 3 is an eddy current image of an impact damaged composite test piece to demonstrate one of the practical uses of the approach. The scan area is 2.5 inch (6.35 cm) by 0.5 inch (1.27 cm). The data acquisition interval is 0.01 inch (0.254 mm) for both x and y axis, i.e. 250 points by 50 points. Although the current application is for eddy current data acquisition and image generation, the approach can be easily applied to ultrasonic imaging and other engineering systems. The only modification necessary is to substitute the eddy current measuring subsystem with the desired engineering instrumentation.

EXPERIMENTS AND RESULTS

The time needed for a given engineering process at a desired location depends only on the process itself. This time is the same regardless of mechanical approaches. The tangible benefit is the decreased scan time in the interrupt-based approach versus the point-to-point approach. The main effect for the improvement is that the interrupt driven scan maintains a constant speed along the scanning axis during data acquisition, while the point-to-point scan must stop at designated intervals. Thus the experiments and data acquisition software have been setup to enable the recording and comparison of the time required for both the interrupt and point-to-point approaches. Identical scan parameters such as scan speeds and index sizes were used for both approaches for the comparisons.

The test specimen used for the bench mark tests is a 3 inch (76.2 mm) by 4 inch (101.6 mm) aluminum block. The scanning velocity for x-axis and y-axis is set to be 0.5 inch (12.7 mm) per second. Three scan configurations are used for the experiments: (1) x-step size of 0.025 inch (0.635 mm) and y-step size of 0.025 inch (0.635 mm); (2) x-step size of 0.025 inch (0.635 mm) and y-step size of 0.050 inch (1.27 mm); and (3) x-step size of 0.050 inch (1.27 mm) and y-step size of 0.050 inch (1.27 mm). Identical scans are performed at least two times to ensure a proper estimate of scan times. The scan time for the described tests are recorded and compared. The average scan time variation is approximately ten seconds or less. The results from the tests are tabulated in Table 1.

The ratio of the two scan times, point-to-point scan time over interrupt scan time, is calculated for each of the three tests conducted. This ratio is used as the measure of the improvement factor. As shown in Table 1 there is approximately a factor of two improvement in scan time for the interrupt scan as compared to the point-to-point scan. This improvement factor is a function of the scan configuration such as scanning speed, data acquisition interval, and specimen size, etc. However, the dominate factor is the number of acquisition points along the scanning axis.



Figure 3. A typical eddy current image of an impact damaged composite test piece (2.5" x 0.5" scan envelop).

Table 1. Comparison of the experimental results from point-to-point and interrupt scans.

Test Configuration Scan Area = 4" x 3" X velocity = 0.5"/sec Y velocity = 0.5"/sec	Point-to-Point	Interrupt	Improvement Factor
x-step size = 0.025" y-step size = 0.025" (160 pts x 120 pts)	2375.49 sec	1013.52 sec	2.34
x-step size = 0.025" y-step size = 0.050" (160 pts x 60 pts)	1192.42 sec	514.01 sec	2.32
x-step size = 0.050" y-step size = 0.050" (80 pts x 60 pts)	1015.27 sec	526.10 sec	1.93

CONCLUSION AND DISCUSSION

In summary, we have (1) proposed a new approach by using an interrupt control mechanism to improve pointwise engineering systems; (2) performed experiments to prove the concept, and (3) verified the practical application aspect by implementing the concept in an eddy current imaging system. Technical improvements include (1) optimized operating parameters; (2) reduced wear and tear on the mechanical system; (3) increased throughput; and (4) improved accuracy for data acquisition and image generation. These improvements translate to increased productivity and reduced cost in engineering operations.

IBM-PCs and their compatibles are gaining in popularity as system controllers and host computers for many mechanical control, instrument control, and signal processing boards. IBM-PC based manufacturing and test/measuring systems thus are routinely being developed, introduced and implemented in various industries. This new approach of using interrupts to initiate and synchronize engineering events has immense commercial potentials; it can be applied to engineering systems in manufacturing, testing, evaluation, and monitoring such as material dispensing, packaging, sorting, and many other industrial applications.

In conclusion, we have demonstrated the use of an interrupt control mechanism for PC-based eddy current NDE data acquisition and image generation. IBM-PCs and their compatibles are gaining in popularity as the system controllers and hosts for mechanical and instrument control boards used in many manufacturing and measuring systems. This new approach of using interrupts to initiate and synchronize engineering events has tremendous commercial potential; it can be applied to systems in manufacturing, evaluation, and monitoring. Specific examples are material dispensing, packaging, sorting, and many other applications.

REFERENCES

1. E. J. Chern, *Materials Evaluation*, Vol. 49, No. 9, September 1991, p1228.
2. M. J. Golis, *An Introduction to Nondestructive Testing*, American Society for Nondestructive Testing, Inc., 1991.
3. E. J. Chern and D. W. Butler, *NASA Tech Briefs*, Vol. 16, No. 9, September 1992, p44.
4. E. J. Chern, to be published in the proceedings of the 1992 *Review of Progress in Quantitative NDE*, Vol. 12, Plenum Press, New York, 1993.

APPENDIX

```

/* Program: INTERRUPT.C */

/*****/
/* ~do_interrupt_scan */
/*****
This is the main scanning routine in which the data will be recorded using interrupts generated on the actual
position of the X-Y scanning table. Global 2-D array, Data_array, will be sized and allocated with the
values defined by the user. If there is not enough available memory to store the data for the entire scan, the
user will be asked to reduce the size of the scan or abort. Array indexing will be done by incrementing the
array pointer whenever a position interrupt occurs. Since these events are controlled as to occur only at
valid data acquisition points within the scan, data for each scan line is stored sequentially in the array. The
data associated with each scan line in the array can be decoded by knowing the total number of data points
for each scan line. This number will vary depending on the size of the scan and the increment size. The
velocity of the scan has to be limited in such a way as to allow the impedance meter enough time to take a
valid reading.
*****/

void do_interrupt_scan (int scan_type)
{
DATA_VAL huge *mem_ptr ;           /* Huge pointer for traversing data array */
char data_buf [MAX_LEN] ;         /* Response string for impedance meter */
char cmd_str [MAX_LEN] ;         /* Command string for PMAC */
int abort = FALSE ;              /* Flag to abort scan */
int done = FALSE ;              /* Flag for normal exit of scan */
int key ;                        /* Key pressed on keyboard */
/* initialize global flags */
Breq_flag = 0 ;                 /* Interrupt flag for PMAC */
Equ1_flag = 0 ;                 /* " " " */
Inpos_flag = 0 ;               /* " " " */
Endofscan = FALSE ;           /* Flag to indicate end of scan */
/* Allocate memory for 2-D data array */

if ((Data_array = define_data_array (Num_Y_pts, Num_X_pts)) == NULL)
{
display_mem_error () ;
exit_program (ERROR) ;
}
mem_ptr = Data_array ;         /* Set pointer to beginning of data array */
/* Send PLC programs and definitions to PMAC */
dnload_pmac_defines () ;
dnload_PLC_0 () ;
dnload_PLC_1 () ;
/* Setup instruments and position the probe */
setup_HP4194 () ;
display_positioning_msg () ;
/* Initialize PMAC and the host PC to accept interrupts */
init_interrupt_mode () ;
/* Begin data acquisition */
send_pmac_cmd ("R") ;         /* Send RUN command to PMAC card */
while (!done)
{
/* Now process the interrupts */
if(Equ1_flag)                 /* At data measurement point */
{
Equ1_flag = 0 ;             /* Reset Interrupt flag */
/* Get reading from impedance meter and store it in memory */
}
}
}
}

```

```

    Receive (BRD, Imp-meter, data_buf, MAX_LEN, STOPend) ;
    mem_ptr->val_1 = atof (data_buf) ;          /* store impedance value */
    mem_ptr++ ;
}
if (Breq_flag)                                /* PMAC ready to receive next command */
{
    Breq_flag = 0 ;                            /* Reset interrupt flag */
    send_next_pmac_move () ;                  /* Send a move sequence to PMAC */
}
if (Inpos_flag)                                /* All motion and move timer stopped */
{
    Inpos_flag = 0 ;                          /* Reset interrupt flag */
    if (Endofscan)                            /* Terminate loop if scan complete */
        done = TRUE ;
}
if (kbhit ())                                  /* Process user requests */
    if (getch() == ESC)                       /* Check for ABORT request */
    {
        send_pmac_cmd ("H") ;                /* Halt the scan */
        done = abort = TRUE ;               /* Terminate loop */
    }
}
/* End of data acquisition loop */
/* Restore the original PC interrupt vectors and store acquired data to disk */
restore_interrupt_vector () ;
if (!abort)
{
    write_data_to_disk (Data_file) ;
    create_header_file (Data_file, scan_type) ;
}
/* Free the memory used for data storage */

farfree ((void *) Data_array) ;
}
/** End of do_interrupt_scan **/

/*****/
/* ~init_interrupt_mode */
/*****/
This routine will setup the Programmable Interrupt Controller (PIC) on both the PMAC card and the IBM-PC. PMAC will be interrupting the IBM_PC on one of the interrupt request lines (IRQ). The IRQ line is defined by PC_IRQ in EC_DEFIN.H. Interrupt request levels 1 (IR1) and 5 (IR5), which correspond to the buffer request and equ1 lines on PMAC's PIC, will be used to generate interrupt pulses that are to be sent to the IBM-PC, on the selected IRQ line. When these interrupt pulses are acknowledged by the PC's PIC, interrupt service routines will be activated to record data and send motion commands to PMAC's rotary buffer.
IMPORTANT: in order for this routine to function properly, jumpers on the PMAC card must be installed to reflect the PC IRQ line that will be used to interrupt the host PC. A jumper at E65 must also be installed to electrically connect the equ1 line to PMAC's PIC.
*****/
void init_interrupt_mode (void)
{
    /* Save original interrupt vector for PC IRQ line, so it can be restored when the program terminates. If this is not done, the default handler for this IRQ will not function without rebooting. */
    disable () ;                                /* disable interrupts until done */
    Old_int_vector = getvect (PC_INT) ;         /* save original interrupt vector */
    setvect (PC_INT, pmac_isr) ;               /* write in new interrupt vector */
}

```

```

/* Save original mask value for the PC PIC's interrupt request register and enable interrupt request level
used by PMAC */
#ifdef PC_PIC_2
Old_int_status = inportb (PC_PIC2_ODD);          /* Save old mask value of PIC #2 */
outportb (PC_PIC2_ODD,(Old_int_status & PC_MASK)); /* Enable IRQ used by
PMAC */
#else
Old_int_status = inportb (PC_PIC1_ODD);          /* Save old mask value PIC #1 */
outportb (PC_PIC1_ODD,(Old_int_status & PC_MASK)); /* Enable IRQ used by PMAC */
#endif
/* Set up PMAC's PIC so it can generate interrupt pulses when the BREQ or EQU1 lines go high. */
outportb (BASE, FLUSH);                          /* Flush PMAC's interrupt control register */
outportb (PMAC_PIC_EVEN, EDGE_TRIG);             /* Set edge triggered mode (ICW1) */
outportb (PMAC_PIC_ODD, BUS_VECTOR);            /* Vector for data bus (ICW2) */
outportb (PMAC_PIC_ODD, MODE_8086);             /* Set up for 8086 mode (ICW4) */
outportb (PMAC_PIC_ODD, PMAC_MASK);            /* Unmask IR1 (BREQ) & IR5 (EQU1) (OCW1) */
outportb (BASE, DSP_READ);                      /* Enable PMAC DSP read */
enable ();                                       /* enable interrupts, done! */
}
/** End of init_interrupt_mode **/

/*****/
/* ~pmac_isr */
/*****/
This is the interrupt service routine for a PMAC generated hardware interrupt in the host. PMAC will be
interrupting the host PC on one of IRQ lines. In order for this routine to function, a jumper must be set on
the PMAC card to indicate which IRQ line is being used, as defined by PC_IRQ in EC_DEFIN.H. A
jumper must also be installed at E65 on the PMAC card to connect the compare-equals signal (EQU1), to
PMAC's PIC (8259A). This routine will test the In Service Register (ISR) of PMAC's PIC to determine
which event has triggered an interrupt, and then set a flag in the host to indicate that event. Currently the
events that are to be tested for are: the buffer request (BREQ), EQU1 signal, and the in position (IPOS)
signal.
*****/

static void interrupt far pmac_isr (void)
{
char isr_val;                                     /* Value read from PIC's in service register */
disable ();                                     /* Prevents other interrupts */
outportb (PMAC_PIC_ACK, PMAC_NOP);             /* Rising edge of 1st INTA pulse */
outportb (PMAC_PIC_EVEN, PMAC_NOP);           /* Trailing edge of 1st INTA pulse */
outportb (PMAC_PIC_ACK, READ_ISR);            /* Setup to read PIC's ISR */
isr_val = inportb (PMAC_PIC_EVEN);            /* Read PIC's ISR value, generates rising edge of
2nd INTA pulse */
if (isr_val & PMAC_BREQ)                       /* If BREQ event, set a flag */
Breq_flag = 1;
if (isr_val & PMAC_EQU1)                       /* If EQU1 event, set a flag */
Equ1_flag = 1;
if (isr_val & PMAC_IPOS)                       /* If in position, set a flag */
Inpos_flag = 1;
outportb (PMAC_PIC_EVEN, PMAC_NOP);           /* Trailing edge of 2nd INTA pulse */
#ifdef PC_PIC_2
outportb (PC_PIC2_EVEN, PC_EOI);              /* Send end of interrupt byte to PC PIC #2 */
#endif
outportb (PC_PIC1_EVEN, PC_EOI);              /* Send end of interrupt byte to PC PIC #1 */
enable ();                                     /* Re-enables other interrupts */
}
/** End of pmac_isr **/

```

```

/*****/
/* ~dnload_PLC_0 */
/*****/

```

This routine will down load a PLC program to PMAC to generate interrupts when the X axis encoder reaches set positions at even increments in both scanning directions. The position-compare-function of PMAC's DSP-Gate array is utilized. Positions at which interrupts are to occur are calculated and then preloaded into the compare register, thus generating an interrupt in the host computer when the encoder value and compare register value are equal.

```

/*****/

```

```

void dnload_PLC_0 (void)
{
send_pmac_cmd ("OPEN PLC 0");           /* Open buffer */
send_pmac_cmd ("CLEAR");               /* Clear it */
send_pmac_cmd ("IF(M116=1)");          /* ENC1 EQU flag bit set? */
send_pmac_cmd ("IF(P1=0)");           /* Negative scan direction ? */
send_pmac_cmd ("WHILE (M101=M105)");   /* Wait for update of position */
send_pmac_cmd ("ENDWHILE");
send_pmac_cmd ("M103=M105");           /* Load next EQU position */
send_pmac_cmd ("M105=M105-P101");      /* Calc. following EQU position */
send_pmac_cmd ("IF(M105=P201-P101+P301)"); /* End of neg. dir. scan line? */
send_pmac_cmd ("M105=M105+P101");      /* Prepare pos. dir. position */
send_pmac_cmd ("P1=1");               /* Set positive direction flag */
send_pmac_cmd ("ENDIF");
send_pmac_cmd ("ELSE");               /* Moving in positive direction */
send_pmac_cmd ("WHILE (M101=M105)");   /* Wait for update of position */
send_pmac_cmd ("ENDWHILE");
send_pmac_cmd ("M103=M105");           /* Load next EQU position */
send_pmac_cmd ("M105=M105+P101");      /* Calc. following EQU position */
send_pmac_cmd ("IF(M105=P301+P101)");  /* End of pos. dir. scan line? */
send_pmac_cmd ("M105=M105-P101");      /* Prepare neg. dir. position */
send_pmac_cmd ("P1=0");               /* Set negative direction flag */
send_pmac_cmd ("ENDIF");
send_pmac_cmd ("ENDIF");
send_pmac_cmd ("M111=0");              /* Clear and set latch control bit */
send_pmac_cmd ("M111=1");              /* to clear latched flag */
send_pmac_cmd ("ENDIF");
send_pmac_cmd ("CLOSE");
while (getline(Response_buf));        /* Clear PMAC's data register */
}
/** End of dnload_PLC_0 **/

```

```

/*****/
/* ~nload_PLC_1 */
/*****/

```

This routine will down load a PLC program to PMAC to activate the DSP-Gate array registers on the PMAC card. The DSP-Gate array will be used in PLC program dnload_PLC_1 and initialize those registers with the proper values to start the interrupt generating sequence. This PLC program is executed only once and then disables itself.

```

/*****/

```

```

void dnload_PLC_1 (void)
{
send_pmac_cmd ("OPEN PLC 1");           /* Open buffer */
send_pmac_cmd ("CLEAR");               /* Clear it */

```

```

send_pmac_cmd ("M111=0");
send_pmac_cmd ("M111=1");
send_pmac_cmd ("M112=1");
send_pmac_cmd ("M113=0");
send_pmac_cmd ("P1=0");
send_pmac_cmd ("P301=M101");
send_pmac_cmd ("M105=M101");
send_pmac_cmd ("M103=M105");
send_pmac_cmd ("M105=M105-P101");
send_pmac_cmd ("ENABLE PLC 0");
send_pmac_cmd ("DISABLE PLC 1");
send_pmac_cmd ("CLOSE");
while (getline(Response_buf));
}

/* Make sure ENC1 EQU flag latch */
/* control bit is reset */
/* Enable EQU output */
/* Set EQU output to high true */
/* Clear direction flag */
/* Get starting position from ENC1 */
/* Init. counter to starting position */
/* Load first EQU position */
/* Calc. following EQU position */

/* Clear PMAC data register */

/** End of dnlod_PLC_1 **/

```

omit

**ADVANCED MATERIALS PART 3:
PLASTICS, POLYMERS, AND RUBBERS**



OMIT

ELECTRO-EXPULSIVE SEPARATION SYSTEM

This paper was withdrawn from presentation

PRECEDING PAGE BLANK NOT FILMED

546-27
150516
P. 8

N 93 - 25607

ELASTOMER COMPOUND DEVELOPED FOR HIGH WEAR APPLICATIONS

D. Crawford, H. Feuer, D. Flanagan, G. Rodriguez, A. Teets, & P. Touchet
Engineering Materials & Coatings Division
Materials Directorate
Army Research Laboratory
Ft. Belvoir, VA 22060

ABSTRACT

The U. S. Army is currently spending 300 million dollars per year replacing rubber track pads. An experimental rubber compound has been developed which exhibits 2 to 3 times greater service life than standard production pad compounds. To improve the service life of the tank track pads various aspects of rubber chemistry were explored including polymer, curing and reinforcing systems. Compounds that exhibited superior physical properties based on laboratory data were then fabricated into tank pads and field tested. This paper will discuss the compounding studies, laboratory data and field testing that led to the high wear elastomer compound.

BACKGROUND

Track laying vehicles, wherein a continuous track is constantly laid down in the direction of movement of the associated vehicle, are well known. Examples of such track laying devices are the military tanks and personnel carriers. Such devices have an endless track with a plurality of linked metal track shoes. These military tracked vehicles are equipped with rubber track pads, rubber blocks or endless-band rubber track to reduce shock, noise, wear and damage to road surfaces. The endless tracks render the vehicles operational in rough, uneven terrain when necessary under military conditions. The vehicle, however, also travels over roads and hard surfaces, therefore, the elastomeric components of the endless track should be of the type that wears well under abrasion and rough terrain.

Historically, field performance of these elastomeric components have been poor, especially for the medium to heavy tonnage tracked vehicles. The problem is further complicated with off-the-road service conditions where pads fail at a much faster rate. The severity of the wear is more pronounced on the M-1 main battle tank than on the older M-60 due to an increase in weight and acceleration, while using a pad with a smaller footprint. This produces higher stresses resulting in higher heat build-up. Therefore, costly and frequent replacement is necessary to keep the tanks operational. The cost of maintaining and replacing track pads is consuming about 25% of the total U.S. Army Operational Maintenance Budget or about 300 million dollars per year.

The service life of these tank pads is affected not only by the terrain conditions but also by the speed, weight of the vehicle and track design. While the operational life of the metal components is approximately 5000 miles for heavy vehicles, the average life of the rubber pads is seldom more than 1500 miles under the best circumstances and is usually less than 550 miles under the severest conditions. During service, the elastomeric pad components are adversely affected in several ways. Most common effects of wear include cuts, tears, heat build-up, flex fatigue, and abrasion.

Conventional track pads, based on styrene-butadiene rubber (SBR), usually fail prematurely in service because of excessive wear, blow-out, which then leads to chunks of material leaving the track pad, and rubber-to-metal bond failure. The widespread use of SBR rubber in tank pad applications is primarily due to the relative low cost of the base SBR polymer along with a U.S. Government policy that requires materials used for military applications to have a domestic source. This policy resulted from the non-availability of natural rubber during World War II. Performance specification MIL-T-11891D was approved in 1984 allowing the use of polymers other than SBR for tank track applications. This option provided

infinite alternatives to the compounder to approach the optimization of specific material properties. Until this point, previous efforts since the mid 1960's had only provided incremental improvements in service life of tank pads of about 5%.

APPROACH

To improve field performance of tank pads, one must first identify those properties that are critical and then optimize them. This is by no means a trivial task. Any rubber compounder would agree that to improve performance of tank pads, properties such as cutting and chipping resistance, tear and tensile strength, crack initiation and growth resistance, abrasion resistance, hysteresis and retention of properties at elevated temperatures would have to be improved. To achieve this tremendous task the U.S. Army sponsored a series of investigations involving industry, government and academia. In 1983 the Rubber and Coated Fabrics Research Group, Belvoir RD&E Center (currently the Engineering Materials and Coatings Division of the Army Research Laboratory) was tasked to conduct a series of studies to improve the service life of tank pads.

Compounding Studies

Comprehensive compounding and processing studies were performed to determine what combinations of formulation ingredients and/or mixing variables affect physical properties. Various polymer systems were evaluated including chloroprene, nitrile (NBR), highly saturated nitrile (HNBR), urethane, natural and synthetic polyisoprene, carboxylated nitrile, polybutadiene, SBR and blends of the above polymers.

Selection of the base polymer is critical to obtain specific characteristics of the final product. For example, natural rubber compounds will most likely exhibit superior resistance to tear and lower hysteresis when compared with polybutadiene compounds which in turn provide excellent flexibility, superior resistance to abrasion, crack growth and heat build-up. Nitriles and neoprenes have high resistance to oils and chemicals.

Various fillers and curing systems were explored to enhance physical properties and minimize reversion. Carbon black and novel non-black fillers were used to improve dispersion, increase toughness and abrasion resistance.

Laboratory Physical Testing

The experimental materials underwent extensive physical testing in the laboratory. Physical tests included tensile strength and tear strength at ambient and elevated temperatures, abrasion, cutting and chipping and dynamic tests such as blow-out and flex crack growth. Compounds that exhibited superior physical properties based on laboratory data were then fabricated into tank pads by Caterpillar Tractor Company. These tank pads were subjected to field testing on the Counter Obstacle Vehicle (COV) and the M-60 tank. Simultaneous to the field testing, samples were taken from the fabricated tank pads for additional laboratory testing. In addition to the experimental formulations, standard production SBR pads were also included in the laboratory and field tests. The laboratory data generated from the fabricated pads is shown in Tables I and II. Table III lists the test methods that were used.

Vehicle Field Testing

Two types of vehicles were employed to carry out the field testing. The COV is an engineering type tracked vehicle that weighs about 72 tons. The other vehicle used was the M-60 battle tank weighing about 45 tons. The testing of the T-107 (COV) pads was performed over a severe course designed to combine all possible operational and terrain factors. The COV testing was conducted at the Engineering Proving Grounds at Ft. Belvoir, Virginia. The testing was performed from November 1986 through April 1987. A

total of 1600 miles was accumulated on the T-107 pads.

The test plan for the M-60 (T-142 pads) field test was designed to include three phases consisting of a 2000 mile paved surface course, a 900 mile hilly cross-country course and a 1000 mile combination course. The M-60 test was conducted at the U.S. Proving Grounds in Yuma, Arizona. The testing began in October 1986 and was completed in May 1988.

RESULTS AND DISCUSSION

NBR-12, an experimental compound based on a highly saturated nitrile elastomer and a novel filler and curing system exhibited superior physical properties based on laboratory data as shown in Tables I and II. The tensile strength of NBR-12 is about 30% higher than that of the commercial pads. The NBR-12 material retained 100% of its original tensile strength after heat aging, compared to about 50% tensile retention for the standard materials. This material also exhibited higher hardness and load bearing capability (see 40% compressibility). SBR rubber with an equivalent hardness is too difficult to mix and process into tank pad configurations. Another significant improvement was achieved on tear resistance. Improvements in tear strength of about 50% were observed at ambient temperature and in some cases tear strength of NBR-12 doubled that of the standard material at elevated temperatures.

Abrasion resistance, a critical tank pad property, was measured by the Tabor and Pico methods. A 24 fold increase in the Tabor abrasion resistance was exhibited by the NBR-12 material. The increase in resistance to abrasion as measured by the Pico method ranged from 300% on COV pads to as much as 600% on the M-60 pads.

A common mode of failure for pads during cross-country operations is cutting and chipping. This property was measured with the Goodrich Cutting and Chipping Machine and both specimen diameter and weight loss were recorded. Test results for both measurements were in excellent agreement and the NBR-12 material provided a 75% improvement for COV pads and over 30% for M-60 pads.

The NBR-12 elastomeric tank pad materials exhibited higher heat build-up but the unique combination of the highly saturated nitrile polymer (HNBR) with the novel curing and reinforcement system produced an unparalleled retention of physical and dynamic properties at high operating temperatures, thus reducing premature failures.

Crack growth resistance relates to the ability of the rubber to deter crack propagation once the rubber has been cut, typical of operation over cross-country terrain. Crack propagation further deteriorates into tear and eventually chunks of rubber can be removed from the pad. Crack growth was measured using a Demattia Flex Tester. The crack growth resistance of NBR-12 showed greater than 400% improvement over the standard COV pad material at ambient temperature while providing at least a 60 fold improvement in crack growth resistance over the standard COV pad material at elevated temperature (250 °F). NBR-12 exhibited a 300% improvement in crack growth resistance at room temperature compared to the standard M-60 pad material and greater than a 30 fold improvement at elevated temperature.

NBR-12 showed exceptional resistance to wear during field testing, extending the service life 2 to 3 times that of standard production pads during the M-60 field test. During the paved course portion of the M-60 field test, standard production pads failed on the average at 1200 miles. The pads with the NBR-12 material were tested to 2000 miles (the maximum duration allowed under the test plan) with out failures. The limited wear of these pads indicated a projected life of 3400 miles, a service level never before achieved by any commercial or other experimental material. Subsequent field tests on M-60 vehicles confirmed the original projections of the service life of the NBR-12 material by exhibiting serviceability beyond 3700 miles on the combination course at Yuma, AZ. Table IV and Figure I, show the improved service life of NBR-12 over the standard production pads.

CONCLUSIONS

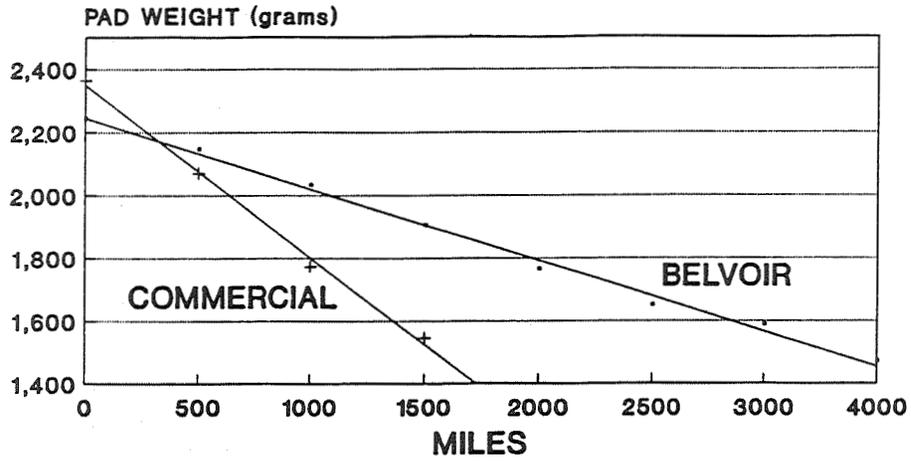
The superior physical properties of NBR-12 have been demonstrated in the laboratory as well as in the field. A U.S. patent has been awarded for this material and patents have been filed in 7 foreign countries. Although initial fabrication cost using NBR-12 is approximately 2.5 times as expensive as standard production pads, a preliminary economic analysis indicates a possible savings on future "higher life" tracks (with 6000 mile life expectancy) currently being designed. Further changes in track design using this improved pad material could contribute significantly to a more reliable and dependable fleet of battle tanks for the U.S. Army while reducing track operating and support costs.

Future research on this high wear elastomer compound is aimed at improving processibility and reducing cost by blending it with other polymers. There is an on-going program to coat NBR-12 and blends of this material onto nylon fabric to be field tested on the Army's Lighter Air Cushion Vehicle (LACV-30). This is a challenging application from the processing standpoint as well as field performance where abrasion resistance and high frequency dynamic flexing or flagellation are critical.

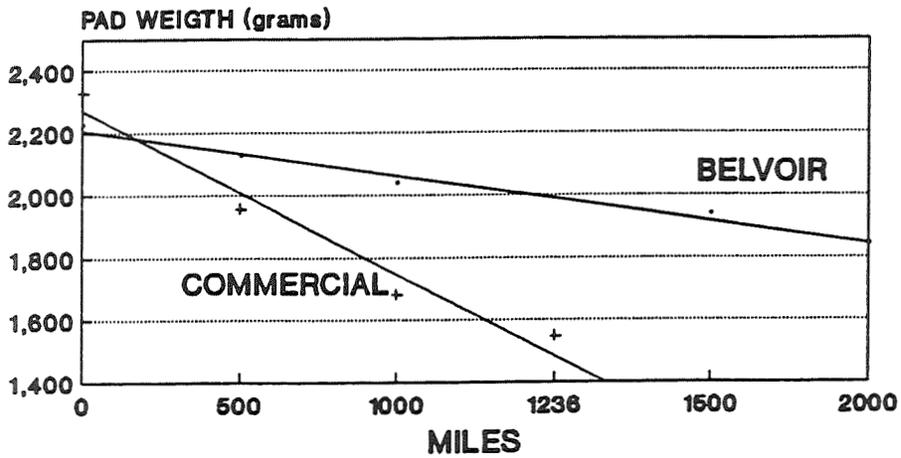
As discussed earlier, NBR-12 is the result of research and development aimed to improve the life of rubber pads used on military vehicles. Applications which require a high resistance to abrasion, resistance to chemicals and fuels, or retention of physical properties at elevated temperatures could benefit from NBR-12's outstanding wear properties. Potential commercial applications for this material include conveyor and "V" belts, treads for off the road tires, gaskets, o-rings and seals for the oil industry, fenders and bumpers on loading docks, shock and vibration pads, and rubber covered rolls for paving equipment.

FIGURE I: Field Performance of T-142 Track Pads

COMBINATION COURSE



PAVED COURSE



HILLY CROSS-COUNTRY

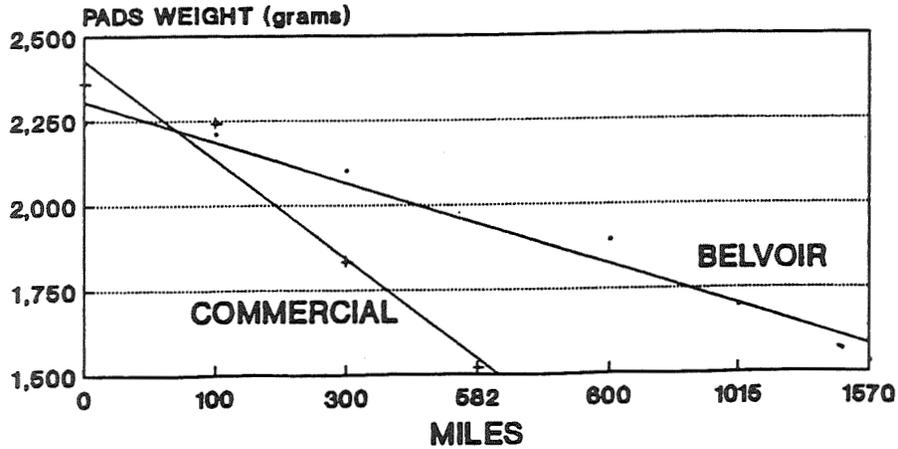


TABLE I: Physical Properties of Materials From T-107 (COV) Pads

MATERIAL I. D. CODE	NBR-12	STD. PAD TP-AF	NAT-150	NN-9
PROPERTIES				
ORIGINAL				
TENSILE STRENGTH, psi	3960	3075	4424	4293
200% MODULUS, psi	841	1102	1944	1566
ELONGATION, %	570	460	380	430
SHORE A HARDNESS, POINTS	78	70	87	84
BASHORE REBOUND, %	30	30	52	51
SPECIFIC GRAVITY	1.1373	1.1626	1.1426	1.1396
40% COMPRESSIBILITY, psi	972	511	1291	1073
TEAR STRENGTH, DIE C				
ROOM TEMP, lb/in	449	268	384	347
AT 250 °F, lb/in	234	141	254	204
OVEN AGED 70 HRS AT 250°F				
ELONGATION RETENTION, %	81	33	71	84
TENSILE RETENTION, %	100	82	59	72
ABRASION				
TABOR, GRAMS/1000 CYCLES	0.0015	0.0434	0.5222	0.4344
PICO RATING	691	179	348	221
GOODRICH CUTTING AND CHIPPING				
DIAMETER LOSS, cm	0.137	0.582	0.516	0.361
WEIGHT LOSS, GRAMS	0.75	3.276	2.965	2.066
DYNAMIC PROPERTIES				
GOODRICH FLEX				
BLOW OUT TIME, MINUTES	5	21	58	120
TEMPERATURE RISE				
INTERNAL, °C	84	72	46	101
EXTERNAL SURFACE, °C	58	28	15	43
DEMATTIA FLEX				
CRACK GROWTH				
UNAGED, in/MINUTE	0.060	0.350	0.300	0.169
20 HR @ 250 °F, in/MIN.	0.118	8.110	0.142	0.382
CRACK INITIATION, 1000 CYCLES	110	10	39	112

NOTES:

- NBR-12 - Experimental pad material based on HNBR polymer.
- STD PAD - Commercial pad material based on SBR rubber.
- TP-150 - Experimental pad material based on natural rubber.
- NN-9 - Experimental pad material based on a nat./HNBR blend.

TABLE II: Physical properties of Materials From T-142 (M-60)Pads

MATERIAL I. D. CODE	NBR-12	STD. PAD TP-A	TP-I	TP-K
PROPERTIES				
ORIGINAL				
TENSILE STRENGTH, psi	3960	2814	3437	4119
200% MODULUS, psi	841	580	1015	1059
ELONGATION, %	570	510	520	500
SHORE A HARDNESS, POINTS	78	69	73	66
BASHORE REBOUND, %	30	36	23	49
SPECIFIC GRAVITY	1.1373	1.1559	1.2314	1.1210
40% COMPRESSIBILITY, psi	972	508	508	450
TEAR STRENGTH, DIE C				
ROOM TEMP, lb/in	449	310	535	554
AT 250 °F, lb/in	234	115	345	320
OVEN AGED 70 HRS AT 250°F				
ELONGATION RETENTION, %	81	38	33	27
TENSILE RETENTION, %	100	57	41	27
ABRASION				
TABOR, GRAMS/1000 CYCLES	0.0015	0.0364	0.0257	0.0338
PICO RATING	691	101	120	131
GOODRICH CUTTING AND CHIPPING				
DIAMETER LOSS, cm	0.137	0.208	0.495	0.610
WEIGHT LOSS, GRAMS	0.75	1.145	2.916	3.228
DYNAMIC PROPERTIES				
GOODRICH FLEX				
BLOW OUT TIME, MINUTES	5	37	37	120
TEMPERATURE RISE				
INTERNAL, °C	84	50	74	37
EXTERNAL SURFACE, °C	58	30	32	15
DEMATTIA FLEX				
CRACK GROWTH				
UNAGED, in/MINUTE	0.060	0.240	0.060	0.118
20 HR @ 250 °F, in/MIN.	0.118	4.331	0.300	0.272
CRACK INITIATION,				
1000 CYCLES	110	18	48	161

NOTES:

- NBR-12 - Experimental pad material based on HNBR polymer.
- STD PAD, TP-A - Commercial pad material based on SBR rubber.
- TP-I - Experimental pad material based on chloroprene rub.
- TP-K - Experimental blend pad material from industry.

TABLE III: Test Methods

PROPERTY TESTED	TEST METHOD ASTM
ORIGINAL PROPERTIES SPECIFIC GRAVITY TENSILE, ELONGATION & MODULUS HARDNESS, SHORE A RESILIENCE, BASHORE REBOUND TEAR STRENGTH, DIE C ABRASION, TABOR ABRASION, PICO	D 792, PARA 15 D 412 D 1415 D 2632 D 624, DIE C D 3389 D 2228
PROPERTIES RUN AT 250 °F TEAR STRENGTH, DIE C TENSILE AND ELONGATION RETENTION	D 573 & D 624 D 573 & D 412
FLEX FATIGUE TESTS DEMATTIA CRACK INITIATION CRACK GROWTH, UNAGED CRACK GROWTH, AGED 20 HR @ 250 °F	D 430, METH B D 813 "
GOODRICH FLEX, BLOW OUT & TEMP. RISE	D 623

TABLE IV: Field Performance of Tank Pads

MATERIAL I D	M-60 PADS TESTED IN YUMA			COV PADS
	PAVED ROAD (Miles)	HILLY X-COUNTRY (Miles)	COMBINATION (Miles)	BELVOIR COURSE (Miles)
NBR-12	3402	1251	3302	2802
STD PADS	1237	530	1351	1305
TP-I	2002	470	1451	----
TP-K	1590	380	1321	----
NAT-150	920	710	1526	2101
NN-9	1237	681	1701	2052

547-35

158517

p-5

N93-25608

DYNAMIC HARDNESS TESTER AND CURE METER

Walter M. Madigosky and Ralph B. Fiorito
Naval Surface Warfare Center, MS R31
10901 New Hampshire Avenue
Silver Spring, MD 20903-5000

ABSTRACT

The Shore hardness tester is used extensively throughout industry to determine the static modulus of materials. The new apparatus described here extends the capability of an indentor type tester into the dynamic regime, and provides a measurement of the dynamic shear or Young's modulus and loss factor as a function of frequency. The instrument, model and data of typical rubber samples are given and compared to other dynamic measurements.

INTRODUCTION

The standard durometer has been used to measure the static "hardness" or resistance of materials to indentation. Various static models have been developed to quantify these measurements by relating the static Young's modulus and Poisson ratio of a viscoelastic material to the force of indentation, the penetration depth, and size (radius) of the indentor.

To our knowledge, an equivalent satisfactory dynamic model in which the indentation force is time dependent has not been developed analytically. Various empirical models, however, have been examined to model the interaction of transducers placed in contact with a viscoelastic slab to excite shear and compressional waves in the material. R. von Gierke¹ e.g. has produced such a model to examine the acoustic properties of living tissue. He assumes the transducer behaves as a hemispherical radiator whose diameter equals that of the cylindrical transducers used in his experiments. In this model, the radiation impedance of a vibrating sphere is related to the complex dynamical shear and compressional moduli of the material. This model, under certain approximations, can be used to infer the dynamic Young's modulus, E' , and loss factor, δ . To perform these measurements the amplitude and phase of the driving force and acceleration of the transducer which is in contact with the material are measured.

We are examining the validity and applicability of such models experimentally by comparing the predicted results for E' and δ to those obtained by more conventional methods such as those developed by Madigosky and Lee², or other type of dynamic mechanical testing apparatus such as the dynamic mechanical thermal analyser (DMTA).

The goal of this study is to develop a model and simple apparatus which will allow a dynamical measurement of E' and δ of a material in a manufacturing environment.

THEORY

The theoretical model we will first examine for measuring the dynamic response of a viscoelastic material to intrusion by a perturbing force is that developed by von Gierke to examine the acoustic properties of living tissue. The theory considers the medium to be homogeneous, isotropic, compressible, and viscoelastic. The probe perturbing the medium is a cylindrical transducer placed in contact with the surface of the medium. The transducer-medium interaction is modeled as a hemispherical radiator whose diameter equals that of the transducer.

The complex impedance of the medium which is defined as the ratio of force over velocity is related analytically to the complex shear and compressional moduli of the material which are defined as:

$$\mu = \mu_1 + j\omega\mu_2 = \mu_1(1 + \delta_S) \quad (1)$$

$$\lambda = \lambda_1 + j\omega\lambda_2 = \lambda_1(1 + \delta_L) \quad (2)$$

where $\omega\mu_2/\mu_1 \equiv \delta_S$ is the shear loss factor and $\omega\lambda_2/\lambda_1 \equiv \delta_L$ is the compressional loss factor.

The Young's Modulus is related to the shear modulus by

$$E \equiv E'(1 + j\delta_S) \equiv 3\mu \quad (3)$$

so that

$$E' \equiv 3\mu_1 \quad (4)$$

The radiation impedance Z_S of the sphere is

$$Z_S = -4\pi\rho\omega a^3 \cdot \left\{ \left(1 - \frac{3j}{ah} - \frac{3}{a^2 h^2}\right) - 2\left(\frac{j}{ah} + \frac{1}{a^2 h^2}\right) \cdot \left[3 - \frac{a^2 k^2}{(jak + 1)}\right] \right\} / \left\{ \left(\frac{1}{a^2 h^2} + \frac{j}{ah}\right) \frac{a^2 k^2}{(jak + 1)} + \left[2 - \frac{a^2 k^2}{(jak + 1)}\right] \right\} \quad (5)$$

where a = radius of the sphere,

$$k^2 = \rho\omega^2 / (2\mu + \lambda) \quad (6)$$

and

$$h^2 = \rho\omega^2 / \mu \quad (7)$$

We have examined several limiting cases of Z_S

a. $\lambda_1 \gg \mu_2, \omega\lambda_2, \lambda_2 \approx 0$

This limit is achieved when $ak \ll 1$, then

$$Z_S \rightarrow Z_2 \equiv -2\pi\rho\omega a^3 \cdot \left(1 - \frac{9j}{ah} - \frac{9}{a^2 h^2}\right) \quad (8)$$

b. high loss materials ($\omega\mu_2 \gg \mu_1$)

$$Z_S \rightarrow Z_3 \equiv -2\pi\rho\omega a^3 \cdot \left(1 + \frac{9(\mu_2/(2\rho\omega a))^{3/2}}{1 + (2\mu_2/\rho\omega a)^{1/2}} - \frac{9j(\mu_2/2\rho\omega a)^{2/2}}{1 + (2\mu_2/\rho\omega a)^{1/2}}\right) \quad (9)$$

c. low loss materials ($\omega\mu_2 \ll \mu_1$)

$$Z_S \rightarrow Z_4 \equiv -2\pi\rho\omega a^3 / 3 \cdot (\omega - 9\mu_1/\rho\omega a^2 - 9j[(\mu_1/\rho a)^2 + \mu_2/a^2])^{1/2} \quad (10)$$

Z_2 and Z_4 will be of primary interest for the materials with properties in the frequency range of interest (50 -1000 Hz).

A comparison of the real part of Z_S , Z_2 and Z_4 as a function of frequency when the sphere radius $a = 1$ mm shows that the formulas for Z_2 and Z_4 are in excellent agreement with Z_S over a wide frequency range, and diverge from Z_S only at very high frequency.

Z_4 since it is straightforward to interpret and use it as the basis for an experimental method to obtain μ_1 and μ_2 .

The impedance Z_4 can be interpreted as that of a simple oscillator. To see this, consider the equation of a forced harmonic oscillator

$$m\ddot{x} + k^*x = F \quad (11)$$

Here, k^* is complex to account for the case of viscous damping. The impedance of the oscillator $Z \equiv F/\dot{x}$ can be obtained by assuming $x = x_0 e^{j\omega t}$ and $F = F_0 e^{j\omega t}$ so that $\dot{x} = j\omega x$ and

$$Z = mj(\omega - k'/\omega m - jk''/\omega m) \quad (12)$$

where $k^* = k' + jk''$

Comparing Z to the equation for Z_4 we see that Z_4 behaves like an oscillator with mass $\equiv m = 2\pi a^3 \rho / 3$ (which is half the mass of a sphere), stiffness $\equiv k' = 6\pi a \mu_1$ and frictional resistance $\equiv k'' = 6\pi a \rho \cdot [(\mu_1/\rho)^2 + \mu_2/a\rho]$ as observed by von Gierke.

If then, a hemispheric impedance head of real mass M is placed in contact with the viscoelastic material, and it is assumed that the impedance of the material is ideally modeled by that of the radiating sphere model under the approximation $Z_S \approx Z_4$ for the frequency range of interest, the acceleration measured by such an impedance head is given by:

$$F/\ddot{x}_{\text{measured}} = M + Z_4/j\omega = (M + m) - k^*/\omega^2 \quad (13)$$

Then, in view of the interpretation of k' and k'' given above, we may write:

$$k' = 6\pi a \mu_1 = \omega^2 (M + m - \text{Re}(F/\ddot{x}_{\text{meas}})) \quad (14)$$

$$k'' = 6\pi a \rho = \omega^2 [(\mu_1/\rho)^2 + \mu_2/a\rho] = \omega^2 m \text{Im}(F/\ddot{x}_{\text{meas}}) \quad (15)$$

We can then use these equations to calculate $E' = 3\mu_1$ and $\delta = \omega\mu_2/\mu_1$.

Similarly the expression for Z_2 can be inverted to produce these values by solving for the variable $1/h$ in equation (8).

EXPERIMENT

We have performed experiments on a variety of viscoelastic materials to test the validity and applicability of the model given above.

To do this we employ a compact Wilcoxon piezoelectric shaker- impedance head to which we can attach indentors of various diameters, a Hewlett-Packard noise generator / spectrum analyzer which drives the shaker with white noise and measures the real and imaginary parts of the acceleration (amplitude and phase of the driving force). With the help of a computer, the complex impedance of the excited material is calculated and the real and imaginary parts are used to produce the Young's modulus and loss factor as described above.

Tables I. and II. show a comparison of experimental results on generic nitrile and urethane rubber samples using the method described above and those obtained using a dynamic mechanical thermal analyser (DMTA) made by Polymer Laboratories. The DMTA was used to measure the Young's modulus and loss factor as a function of temperature using the frequencies 0.3, 1, 3, 10 and 30 Hz. The data was then shifted to the frequency domain by using the standard Williams-Landel-Ferry (WLF) technique.

Table I. Comparison of Results on Nitrile Rubber

F(Hz)	Dynamic Durometer Method		DMTA Method	
	E (m Pa)	Loss Factor	E (m Pa)	Loss Factor
50	4.3	0.3	7.5	0.2
100	5.3	0.47	7.7	0.25
200	7.0	0.34	7.9	0.3
300	7.5	0.31	8.5	0.33
400	7.8	0.25	8.9	0.35
500	8.1	0.27	9.5	0.4

Table II. Comparison of Results on Urethane Rubber

F (Hz)	Dynamic Durometer Method		DMTA Method	
	E (m Pa)	Loss Factor	E (m Pa)	Loss Factor
200	54	0.15	63	0.085
300	59	0.16	66	0.090
400	61	0.10	68	0.094
500	62	0.11	70	0.097
600	70	0.17	74	0.100

These results indicate the good agreement between the two types of measurement. Further studies are being done to test the effect of indenter surface area and the effect of static contact pressure of the indenter.

In addition to the dynamic modulus measurements on cured materials, a study of the change in dynamic modulus as a function of time was made on a two component polyurethane (Techthane 13, Seaward International) as it changed from a viscous mixture to a cured solid. The results are given in Table III. A small disc, 0.8 cm diameter, was used. As can be seen from the data, this technique may be used to accurately determine the rate of cure and the state of cure in materials.

Table III. Dynamic Cure Test of Polyurethane

Time (hours)	Young's Modulus (kPa)	Loss Factor
1	4	1.7
2	51	0.9
3	320	0.68
5	610	0.52
6	770	0.42
7	850	0.41
19	1150	0.38
31	1620	0.37
43	1850	0.37
53	1890	0.36
94	1980	0.34
182 (8 days)	2150	0.33

Similar rate and state of cure measurements can be obtained on epoxies and vulcanized rubbers, thus replacing current methods which provide only relative data.

REFERENCES

1. H. von Gierke et. al., Physics of Vibrations in Living Tissues, J.Appl. Physiology, 4, p. 886-900, 1952.
2. W. Madigosky and G. Lee, Improved Resonance Technique for Materials Characterization, J. Acoust. Soc. Am. 73, p. 1374-77, 1985.

**INSTRUMENTATION FOR MEASUREMENT OF GAS PERMEABILITY
OF POLYMERIC MEMBRANES***

Billy T. Upchurch
NASA Langley Research Center
Hampton, VA 23681

George M. Wood
NASA Langley Research Center
Hampton, VA 23681

Kenneth G. Brown
Old Dominion University
Norfolk, VA 23529

Karen S. Burns
Old Dominion University
Norfolk, VA 23529

ABSTRACT

A mass spectrometric "Dynamic Delta" method for the measurement of gas permeability of polymeric membranes has been developed. The method is universally applicable for measurement of the permeability of any gas through polymeric membrane materials. The usual large sample size of more than 100 square centimeters required for other methods is not necessary for this new method which requires a size less than one square centimeter. The new method should fulfill requirements and find applicability for industrial materials such as food packaging, contact lenses and other commercial materials where gas permeability or permselectivity properties are important.

INTRODUCTION

The gas permeability of a polymeric material is an important physical property which helps determine whether or not a given material might be appropriate for a particular application. Some application examples are food packaging, beverage containers, gas separation processes and oxygen permeability in contact lens materials. The permeability depends upon two physical processes, the solubility of the gas in the material and the diffusivity (D) of the gas through the material. The solubility often follows Henry's Law behavior and can be represented by the Henry's Law constant k. If the thickness of the membrane (L) is taken into account then we can define the permeance or transmissibility by equation 1:

$$P = \frac{Dk}{L} \quad (1)$$

where P is the gas permeation rate through a known area and thickness of material per unit time.

Several analytical techniques have been applied to the problem of determining gas permeability. A few

*This research is a part of the Gas Permeable Polymeric Materials (GPPM) flight experiment, and is being carried out under the auspices of the Office of Commercial Programs, NASA Headquarters, Washington, DC. The GPPM experiment is manifested as part of the SPACEHAB payload on the STS-57 mission in April, 1993, and includes materials from both NASA and an Industrial Guest Investigator, Paragon Optical Co., Mesa, AZ.

are gas specific, such as polarography for oxygen determination¹, the application of platinum and silver electrodes for measurement of oxygen in living tissues¹, and the application of electron spin resonance (ESR) to oxygen permeation in polyethylene². Additional techniques involve the use of spectrophotometry but these involve the incorporation of photosensitive dyes into the materials^{3,4,5}. The American Society of Testing and Materials (ASTM) has adopted three standard test methods; a) a manometric technique, ASTM:D 1434-82⁶, b) a volumetric method, ASTM:D 1434-82⁶, and c) a coulometric method, ASTM:D 3985-81⁶. Each of these standard techniques requires large sample test areas (50 - 100 cm²) and is not capable of differentiating between gases when gas mixtures are being measured.

This paper describes a mass spectrometric technique that minimizes sample area and allows the determination of permeation rates for individual gases in gas mixtures. This technique, termed the "Dynamic Delta" method⁷, is a significant improvement upon standard mass spectrometric techniques. The technique is applicable for any gas and is easy to perform with good reproducibility. We shall demonstrate its appropriateness to be included as a standard industrial test for gas permeability.

METHOD

The "dynamic delta" measurement, is not specific for any particular pressure measuring device or mass spectrometer as long as the response is linear over the range of interest. The basic setup is illustrated in figure 1 with a standard X-Y plotter as the recording device. With this particular configuration we can simultaneously monitor the ion signal as a function of the inlet pressure. The ion signal I_A is proportional to the partial pressure P_A through an experimentally determined proportionality constant S_A . The magnitude of I_A is determined by the sum of the partial pressure of gaseous constituent A introduced into the ion source, the background pressure of A which may already exist in the ion source, and the contribution from an unresolved ion equivalent in mass to that of constituent A.

In general practice, the background signal is experimentally determined prior to introducing a sample into the mass spectrometer. This signal is then subtracted from the total signal prior to calculating the concentration. When the concentration of the sample is large with respect to the background, the quantitative analysis can be carried out with high accuracy. When analyzing at trace levels however, the signal from the sample may approximate that from the background and the accuracy is significantly reduced.

The effects of background on accuracy can be reduced if the measurements are carried out with several samples introduced into the ion source at incrementally increasing pressures and the concentration calculated from the differences in the resulting signals. Since the ion source background signal should not change, its effect will be reduced along with a reduced requirement for making highly accurate measurement of total pressure. Studies in the NASA Langley Instrument Research Division Mass Spectrometry Laboratory determined that the accuracy would increase with the number of difference calculations. If the ion source pressure were allowed to increase in a linear manner and the ion current from a single mass is simultaneously measured, this would be equivalent to obtaining an essentially infinite number of difference measurements and the concentration could be directly determined from a determination of the slope of the resulting line.

When the ion signal, representative of the gas species of interest, is provided as input to the Y axis of the recorder and the inlet pressure is provided as a signal to the X axis the resultant plot should be linear with any deviations from linearity indicating either experimental error or additional phenomena taking place elsewhere in the system. Typical experimental data, using oxygen as an illustration with argon as an internal reference gas are shown in figure 2. These data may be used to calculate the concentration of a particular species from the following relationship:

$$Vol \% = \frac{100 \cdot \tan\theta_1 \cdot \Omega_2 \cdot \mu_1 \cdot \alpha}{\tan\theta_2 \cdot \Omega_1 \cdot \mu_2} \quad (2)$$

where:

- $\tan\theta_1$ = the slope of the line for the analyte gas
- $\tan\theta_2$ = the slope of the line for the reference gas
- Ω_1 = the ion current amplifier resistance used for the reference gas
- Ω_2 = the ion current amplifier resistance used for the analyte gas
- μ_1 = the recorder attenuation used for the reference gas
- μ_2 = the recorder attenuation used for the analyte gas
- α = the relative gas sensitivity for the analyte and reference gases: identical gases $\alpha = 1$.

In the above expression everything but the angle measurements represents properties of the system and equation 2 may be rewritten as:

$$Vol \% = mol\% = 100 \cdot \frac{\tan\theta_1}{\tan\theta_2} \cdot K \quad (3)$$

where K is the experimentally determined proportionality constant for the analysis.

EXPERIMENTAL

Specimen holders, shown in figure 3, were designed to enable the study of both flat and curved surface specimens and were fabricated from brass. Samples were held tightly in place using an O-ring seal to prevent leakage around the polymer specimen. As shown in figure 4, the analysis system was designed to assure identical gas pressure and flow rates on both sides of the membrane. In this design the permeating analyte gas was allowed to flow by the upstream side of the membrane while the carrier gas, in this case argon, flowed by the membrane on the mass spectrometer side. Permeant gas then enters the argon gas stream and is carried to the capillary inlet of the mass spectrometer for analysis as discussed above.

All of the mass spectrometric measurements were made on a 180°, 12.7 cm radius magnetic sector mass spectrometer of the Dempster type. The instrument is equipped with a capillary inlet designed to permit continuous sampling from atmospheric pressure and to provide control of the pressure in the ion source. A representative sample of the gas is introduced directly into the ion source through a gold foil molecular leak and quartz tube. The entire assembly can be maintained at elevated temperature to prevent condensation of volatile gaseous constituents. The exit slit was adjusted to provide a flat top peak.

The gases used were argon, nitrogen, carbon dioxide with stated purities of 99.998% and research grade oxygen. The gravimetric primary standard calibration gas mixture was certified with the following concentrations: 0.4609% O₂; 0.5122% N₂; 0.8077% CO₂ and 98.23% Ar.

The standard reference material was purchased from the National Bureau of Standards (NBS). Fifteen sheets of a poly(ethyleneterephthalate) film were supplied. Three samples, 1.35 cm in diameter, were cut for the mass spectrometry measurements. Contact lens samples were provided by Paragon[®] Optical with a concave radius of 8.00 mm and a chord diameter of 10.00 mm. The center thickness was nominally 0.20 mm and was measured to ± 0.002 mm.

RESULTS AND DISCUSSION

This method is very amenable to computerization utilizing one channel of an A-D converter for the pressure signal and another channel for the ion signal. Computer programs have been written to process data of this type with highly precise results. Raw data is shown in table I to illustrate measurement precision obtained from duplicate runs on two samples of one contact lens material.

Table I: Raw data for oxygen from the Dynamic Delta method for contact lens specimen 5 at 35 °C.

Run	Sample 5a (L = 0.0186cm)	Sample 5b (L = 0.0192 cm)
Peak (ratio)	Slope or Slope Ratio	Slope or Slope Ratio
³⁶ Ar reference	0.7236 ± 0.0039	0.7257 ± 0.0202
³² O impurity	0.1377 ± 0.0099	0.1347 ± 0.0141
³² O/ ³⁶ Ar impurity	0.1903 ± 0.0137	0.1856 ± 0.0201
³⁶ Ar reference	0.7318 ± 0.0019	0.7321 ± 0.0083
³² O total	0.4677 ± 0.0052	0.4380 ± 0.0066
³² O/ ³⁶ Ar total	0.6391 ± 0.0073	0.5983 ± 0.0115
³² O/ ³⁶ Ar corrected	0.4488 ± 0.0155	0.4127 ± 0.0232

The results of the mass spectrometric analyses of some general polymer materials are presented in Table II. The measured value for the SRM 1470 represents the current detection limit of the system using this grade of Ar gas and the small sample size. This detection limit is due to the presence of oxygen in the argon carrier gas. This level of impurity was a problem in trying to determine the permeability of the standard reference material (SRM 1470). We determined the permeability at 23 °C to be $4.52 \pm 0.01 \times 10^{-12}$ mL O₂/cm s mm Hg, which corresponds to the level of oxygen in the argon gas, and is an order of magnitude greater than the NBS certified value of 2.67×10^{-13} mL O₂/cm s mm Hg. It should be noted that the SRM is a barrier material meant as a standard for the testing of barrier materials.

The SRM 1470 is a highly crystalline film, biaxially cold-drawn when it is made, and has a glass transition temperature T_g close to 90 °C or well above room temperature.⁸ These first two characteristics lead to constraints of the polymer chains to reorient.^{9,10} When a glassy polymer is below its glass transition, it is not in a state of true thermodynamic equilibrium, and the permeability and solubility coefficients are more dependent on gas pressure or concentration in polymers and on temperature.¹⁰ These types of materials generally exhibit dual-mode sorption behavior unless their "excess" free or void volume below T_g is small.¹⁰ Dual mode or type II isotherm can result when these glassy polymers absorb gases into pre-existing voids and then behaves as a true solution. With the exception of the silicone membrane, the other materials in Table II are known to have crystalline or semi-crystalline forms and a T_g well above room temperature. Subsequent literature review did not find oxygen permeability data on these materials for comparison. The oxygen permeability of the silicone membrane data is comparable to the product literature¹¹ and other data¹².

The gas permeabilities for the contact lens materials are summarized in tables III and IV for the temperatures 22 °C and 35 °C respectively. With the exception of the two materials of lowest permeability, 6 and 8, the permeabilities at 35 °C for oxygen, nitrogen, and carbon dioxide were significantly greater than permeation rates determined at room temperature.

Our measured carbon dioxide permeation rates are much greater than oxygen or nitrogen permeation rates which is consistent with trends reported in the literature. On average CO₂ permeabilities were five times greater than oxygen permeabilities and ten times greater than nitrogen permeabilities. Carbon dioxide permeation rates are least affected by increased temperature. The Nitrogen and Oxygen permeabilities for materials 2 and 3 are affected the most by increased temperature. At 35 °C these specimens have over a two fold increase in permeation rates. The remaining materials have nitrogen permeation rates only about 50% greater at 35 °C, while the oxygen permeation rates for the corresponding materials have a somewhat smaller increase. This can be explained by the fact that the Henry's law solubility constant of gases in liquids decreases with increasing temperature while the diffusion constant (D) increases with increasing temperature. Consequently, those materials affected least by temperature increases have lower temperature coefficients for a specific gas permeation rate.

TABLE II: Oxygen Permeability Coefficients for Selected Polymer Materials as Measured by Mass Spectrometry

Material	Dk ^a /10 ⁻¹²	T(°C)	Measured ^b P	Published ^b P
Polyamide (Kevlar)	2.41 ± 0.02	19	0.225	NA
Silicone	4410 ± 10	21	410	458 ^c 450-461 ^d
Polysulfone	43.8 ± 0.20	22	405	
Polyester 7d (SRM 1470)	4.52 ± 0.01	23	0.190	0.0113
7a	3.35 ± 0.80	23	0.141	0.0113
7b	3.37 ± 0.73	23	0.142	0.0113
7c	4.05 ± 0.73	23	0.170	0.0113
Polyimide ^e (ODPA/p-PDA) 2% offset	7.82 ± 0.20	22	0.724	N/A
ODPA-p-PDA	8.10 ± 0.20	22	0.750	N/A

a. The units are (cm²/s)(mL O₂/cm²·mmHg)

b. Correction to STP assumed pressure differences from one atmosphere or 760 mm Hg were negligible. The temperature was corrected using T₀ = 273 K and T_{measured(K)} by multiplying Dk by the factor T₀/T_{measured(K)}. The units are Barrer units mL O₂(STP)/cm²·s·cm Hg

c. From GE product literature after converting from room temperature of 25 °C to STP.

d. Rogers, C.E., In Polymer Permeability, Comyn, J., Ed.; Elsevier Applied Science Publishers, London, 1985; Chapter 2.

e. ODPA/p-PDA are acronyms for oxydiphthalic anhydride/ p-phenylenediamine.

TABLE III: Gas Permeability Coefficients^a for Contact Lens Specimens at 22 °C

Material	Oxygen	Nitrogen	Carbon Dioxide
1	15.4 ± 1.4	9.1 ± 1.6	89 ± 4
2	29.7 ± 1.8	9.2 ± 3.3	265 ± 4
3	43.5 ± 1.3	14.4 ± 1.9	357 ± 11
4	62.5 ± 1.8	34.2 ± 1.7	488 ± 13
5	36.5 ± 0.9	20.4 ± 0.9	227 ± 7
6	7.8 ± 0.8	10.2 ± 1.1	13.5 ± 5.0
8	10.4 ± 1.2	10.1 ± 1.1	5.2 ± 1.0
9	64.1 ± 7.1	27.3 ± 2.1	390 ± 9

TABLE IV: Gas Permeability Coefficients^a for Contact Lens Specimens at 35 °C

Material	Oxygen	Nitrogen	Carbon Dioxide
1	29.7 ± 3.6	13.2 ± 3.1	126 ± 5
2	70.2 ± 2.4	32.4 ± 2.4	297 ± 6
3	78.6 ± 2.8	34.4 ± 4.0	408 ± 8
4	99 ± 1	50.7 ± 4.0	555 ± 3
5	46.0 ± 1.5	27.2 ± 0.9	253 ± 7
6	10.7 ± 0.8	23.7 ± 2.4	18.2 ± 2.9
8	6.0 ± 2.1	11.4 ± 1.6	5.6 ± 1.3
9	90 ± 3.5	44.1 ± 1.6	435 ± 9

a. The units are $(\text{cm}^2/\text{s}) \times (\text{mL O}_2/\text{cm}^3 \cdot \text{mmHg}) \times 10^{-11}$

CONCLUSIONS

We believe that the method reported herein is the most accurate method for determining true or intrinsic gaseous permeabilities of any polymeric system. The "Dynamic Delta" method offers significant advantages over standard mass spectrometric techniques including, but not limited to, speed of measurement. By the design of the test cell that we have employed the total pressure gradient is eliminated. The technique permits the use of a very small sample size, 0.5 cm^2 , is selective for specific gases or isotopes, and is sensitivity limited only by the impurity of the carrier gas and the sealing ability of the gasket material. We have been able to demonstrate a leak-free system by using a metal blank. It was determined that flat Teflon gaskets on both

sides of the flat membrane specimens sealed better than the Buna-N or Viton materials. The Buna-N O-ring, however, has the better sealing ability for the contact lens samples.

The ASTM methods have traditionally been done with large sample sizes in order to provide test areas of 50 - 100 cm². Such large sizes meant that the sample had to be placed upon a support. As a result the gas of interest had to permeate through both the sample and the support. In the "Dynamic Delta" method described in this paper the sample size is less than one square centimeter eliminating the need for the support. The permeability of the material can be measured directly. The new method should fulfill requirements and find applicability for industrial materials such as food packaging, contact lenses and other commercial materials where gas permeability or permselectivity properties are important.

REFERENCES

1. Clark, L. C., in Polarographic Oxygen Sensors, ed. Fatt, I, CRC Press, Cleveland, Ohio 1976.
2. Hori, Y, Shimada, S., and Kashiwabara, H., *ESR Studies on Oxidation Processes in Irradiated Polyethylene: I. Diffusion of Oxygen into Amorphous Parts at Low Temperatures*, Polymer Vol. 18, 151-154 (1977).
3. Shaw, G., *Quenching by Oxygen Diffusion of Phosphorescence Emission of Aromatic Molecules in Polymethylmethacrylate*, Trans. Faraday Soc., Vol. 63, 2181-2189, (1967).
4. Petrak, K., *Permeability of Oxygen Through Polymers. I. A Novel Spectrophotocatalytic Method*, J. of Appl. Polym. Sci., Vol. 23, 2365-2371, (1979).
5. Rooney, M. L., and Holland, R. V., *Measuring Oxygen Permeability of Polymer Films by a New Singlet Oxygen Technique*, Angew. Makromol. Chem., Vol. 88, 209-221, (1980).
6. Annual Book of ASTM Standards, 15.09, ASTM: Philadelphia, 1992.
7. Hughes, D. B. and Nowlin, T. D., *Trace Gas Analysis*, presented at the 19th Annual Conference on Mass Spectrometry and Allied Topics, Atlanta, Georgia, 1971.
8. *Standard Reference Materials: SRM 1470; Polyester Film for Oxygen Gas Transmission Measurements*, NBS Spec. Publ., (U.S.) No. 260-58 June, 1979.
9. Stern, S. A., and Trohalaki, S., in Barrier Polymers and Structures, Koros, W. J., ed., American Chemical Society Symposium Series No. 423, Washington, D. C., Chapter 2, (1990).
10. Weinkauff, D. H., and Paul, D. R., in Barrier Polymers and Structures, Koros, W. J., ed., American Chemical Society Symposium Series No. 423, Washington, D. C., Chapter 3, (1990).
11. General Electric Corporation Product Literature, *Permselective Membranes*, Medical Development Operation, Chemical and Medical Division, Membrane Products, Schenectady, N. Y.
12. Rogers, C. E., in Polymer Permeability, Comyn, J. ed., Elsevier Applied Science, London, chapter 2, (1985).

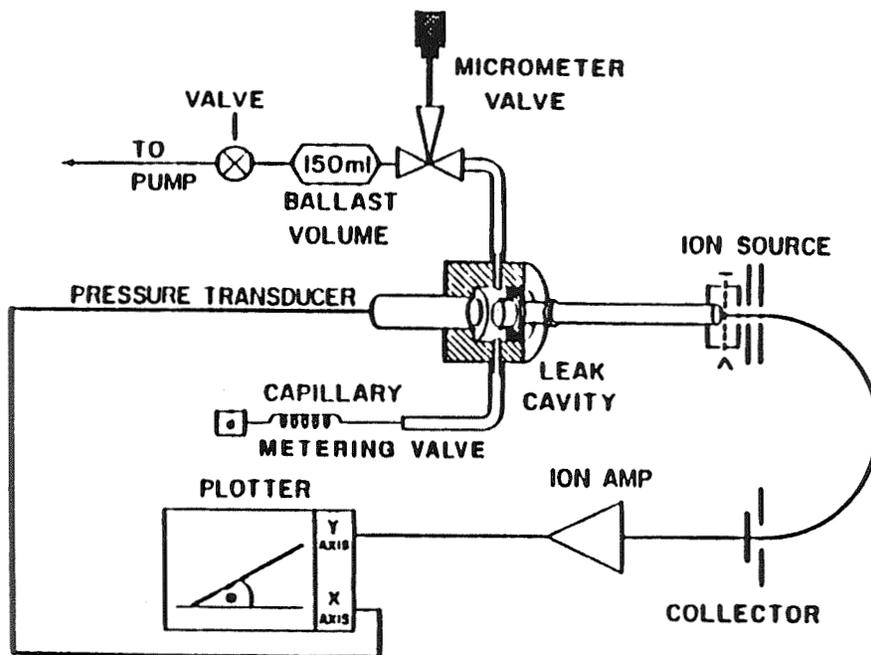


Figure 1. Schematic diagram of mass spectrometer analysis system.

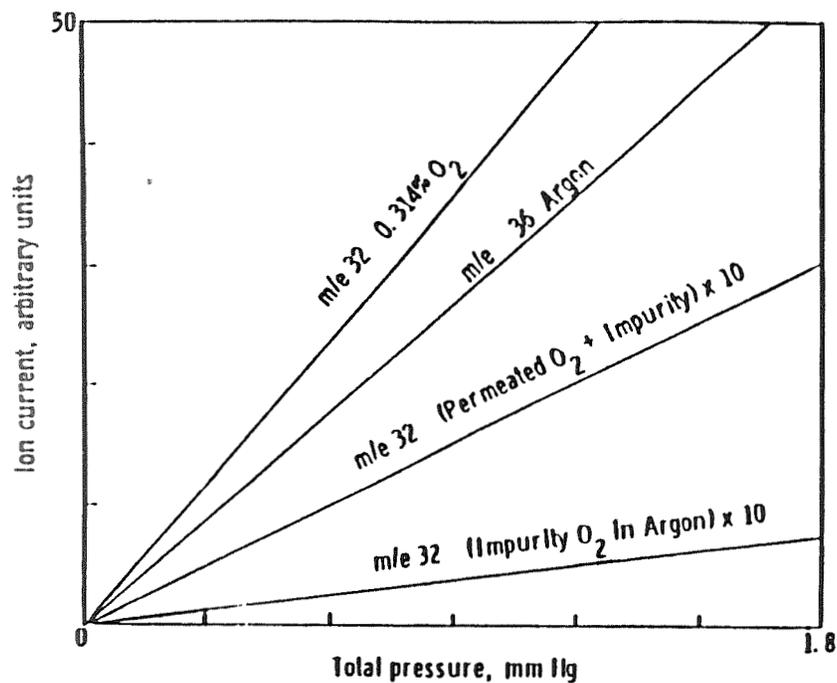


Figure 2. Typical "Dynamic Delta" for oxygen permeation system.

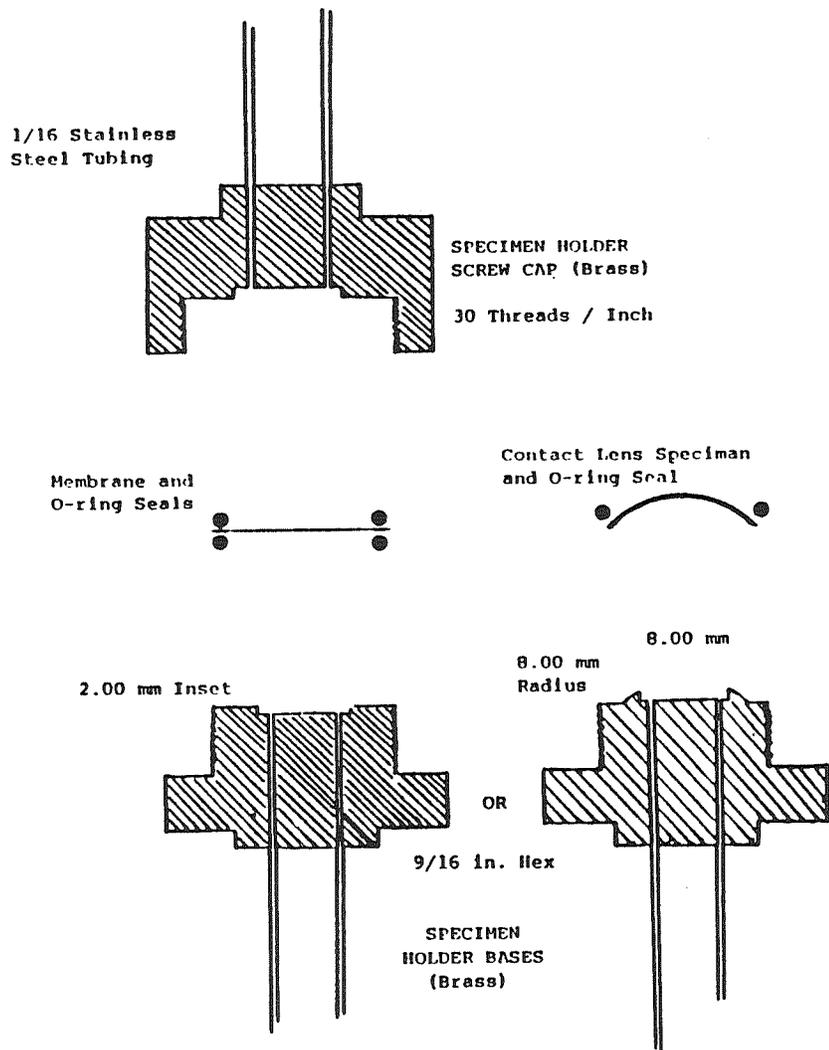


Figure 3. Mass spectrometer sample holder accommodating flat Specimens and contact lens specimens.

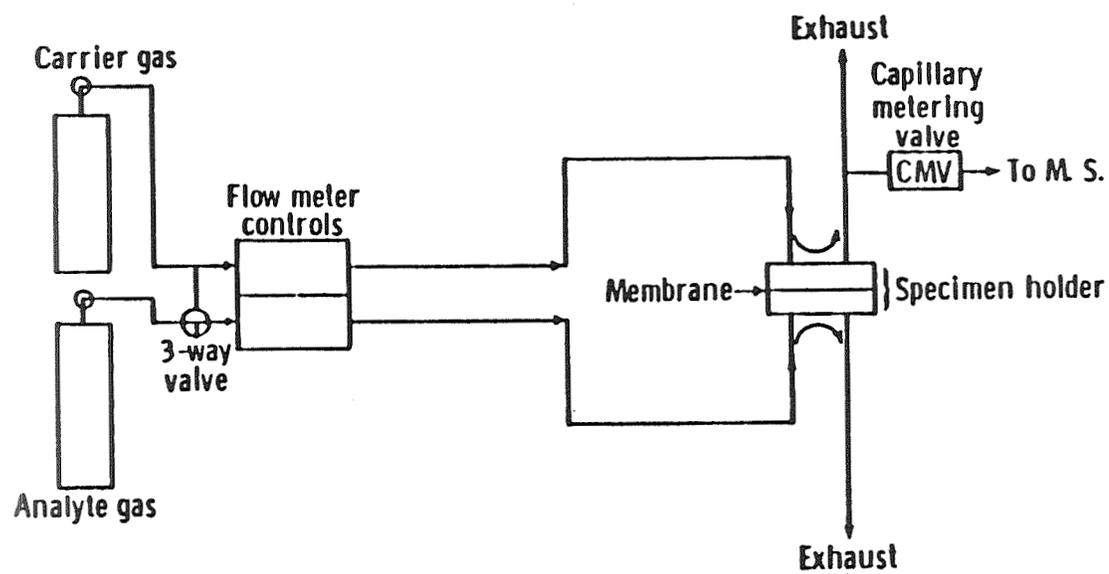


Figure 4. Gas flow schematic for mass spectrometric permeation measurement.

CMIT

**ARTIFICIAL INTELLIGENCE
PART 2**

A SOFTWARE PACKAGE FOR NEURAL NETWORK APPLICATIONS DEVELOPMENT

Robert H. Baran
 Naval Surface Warfare Center
 White Oak (code N51)
 Silver Spring, MD 20903-5000

549-63
 150519
 P-10

ABSTRACT

Original Backprop (Version 1.2) is an MS-DOS package of four stand-alone C-language programs that enable users to develop neural network solutions to a variety of practical problems. *Original Backprop* generates three-layer, feed-forward (series-coupled) networks which map fixed-length input vectors into fixed-length output vectors through an intermediate ("hidden") layer of binary threshold units. Version 1.2 can handle up to 200 input vectors at a time, each having up to 128 real-valued components. The first subprogram, TSET, appends a number (up to 16) of classification bits to each input, thus creating a *training set* of input-output pairs. The second subprogram, BACKPROP, creates a trilayer network to do the prescribed mapping and modifies the weights of its connections incrementally until the training set is learned. The learning algorithm is the "back-propagating error correction procedure" first described by F. Rosenblatt in 1961. The third subprogram, VIEWNET, lets the trained network be examined, tested, and "pruned" (by the deletion of unnecessary hidden units). The fourth subprogram, DONET, makes a TSR routine by which the finished product of the neural net design-and-training exercise can be consulted under other MS-DOS applications.

INTRODUCTION

Recent advances in the manufacture of integrated circuits have led to parallel computers with thousands of microprocessors in a single system and given rise to growing interest in computational methods that support massive parallelism in a natural way. Neural computing aims at (ultimately) achieving human-like performance in computer systems by developing an analogy to the structure and operation of the central nervous system. Computations are executed by simple neuron-like processing units which are interconnected by "synaptic" links of variable strength (or weight). The resurgence of interest in these "connectionist" models, since about 1982, has been said to reflect the total inadequacy of algorithm-driven computing and symbolic artificial intelligence (AI) approaches in dealing with real world problems, speech processing and machine vision being the most cited examples.

On the other hand, neural networks have found abundant use in recent years as practical decision aids which can learn by example to give correct responses to given inputs in situations where, in principle, a set of logical rules could be used to infer the correct response but where, in practice, such rules are difficult to elucidate. Successful case studies have been reported in sonar echo classification, concealed explosives detection, mortgage risk evaluation, medical diagnosis, and many other applications. Much of the popularity surrounding neural network classifiers is a consequence of the relative ease with which they can be trained to substitute for optimally designed expert systems. In most of the interesting applications so far, the powerful new result which enables the neural network to capture the significant factual associations presented in the training data is a learning algorithm called *back-propagation*.

Original Backprop (Version 1.2) is an MS-DOS software package for setting up and training three-layer, feed-forward neural networks to classify input patterns consisting of binary- and real-valued components. It uses the oldest (and least widely known) form of the back-propagation learning algorithm to achieve a degree of flexibility and performance which rivals some of the more popular software-only neural net development products on the market today. Its predecessor (Version 1.1), which was distributed as shareware to members of the international neural networks research community and a tri-services working group, has been applied with some success to the automation of medical diagnosis [1], to active sonar target classification [2], and to

personnel screening. It was also used (to no apparent advantage) in financial forecasting [3]. At the present time, the author is exploring the application of *Original Backprop* to the interpretation of questionnaire data produced by a pre-prototype software package for the prevention and remediation of sexual harassment.

The next section, which explains the historical origins of the learning algorithm, will clarify some points of functionality and terminology which have to be understood before the package can be used effectively. The third section walks the reader through an example problem (using Version 1.1) in which a network with random initial weights is set up and then trained to classify the elements of a training set.

BACKGROUND

The neural network technology of today is based largely on the neuroscience of the 1940s. The classical neuron integrates the pre-synaptic activity of all the neurons influencing it, sending information in the form of electrical impulses down the one-way path formed by its axon. McCulloch and Pitts, in 1943, simplified the neuron to an on/off device, either firing impulses at its peak rate or resting quietly. In 1948, D.O. Hebb theorized that the synaptic weights are modified by a reinforcement control procedure; and he argued that synaptic modification constituted the microscopic, physiological basis of adaptation, learning, and behavioral organization. How could this theory be tested?

By 1954, digital computers had been brought to bear on the problems of brain modeling. The computer was indispensable, because the mathematics involved large numbers of variables and their nonlinear interactions. If "intelligent" behavior was going to emerge from networks of McCulloch-Pitts neurons with Hebbian synapses, it would have to be discovered by computer simulation. No one carried this idea so far, so fast as Frank Rosenblatt, whose discoveries were summarized in a 1961 Cornell Aeronautical Laboratories technical report, *Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms* [4]. In 1962, *Neurodynamics* was published and distributed by Spartan Books (now defunct [5]). The first 300 pages of Rosenblatt's report were devoted to three-layer, series-coupled perceptrons composed of binary threshold units (or McCulloch-Pitts neurons). The front end of the perceptron consists of sensory (S) units on which a pattern of binary digits is impressed. The back end consists of response (R) units which register the classification. Between these is a layer of association (A) units each of which forms a weighted sum of the pattern components and then turns ON (or OFF) when the sum exceeds (does not exceed) a threshold. Similarly, the R-units turn ON or OFF according to the values of the weighted sums that they compute after scanning the A-layer. In Figure 1, some newer terminology is superimposed on this 30-year-old design.

Rosenblatt trained simple perceptrons to solve problems in pattern recognition. For example, let the S-units form a grid-like retina on which horizontal and vertical bars are impressed by turning ON the units in a particular row or column. Let there be only one R-unit which we want to turn ON in response to vertical bars only. The weights of the front-end connections (from S to A) and the back-end connections (from A to R) are initially just random numbers; and the perceptron's initial performance might be correct about half the time. The training process follows a sequence of cycles. A pattern is presented at the front end and propagated through the A-layer to the R-unit. If the response is correct, go on to the next pattern. If incorrect, then change the weights of all A-to-R connections which contribute to the error. The weight change will be negative when the A-unit helps to turn ON the R-unit in response to a horizontal bar, positive when it inhibits the R-unit's response to a vertical bar. As this procedure is repeated again and again, the perceptron's incorrect responses become fewer and fewer.

In 1969, MIT computer scientists Marvin Minsky and Seymour Papert published a book, *Perceptrons*, which is widely regarded as having had a chilling effect on the subject. In the third (1988) edition, Minsky recalls that perceptron research had already reached a dead end [6]. After Rosenblatt died in a boating accident on the Chesapeake Bay, in 1971, Minsky and Papert dedicated the second edition of *Perceptrons* in his memory. Yet the memory of what Rosenblatt accomplished faded quickly in the years that followed as students increasingly accepted the Minsky-Papert perceptron as a substitute for the original--and found it lacking in problem-solving ability.

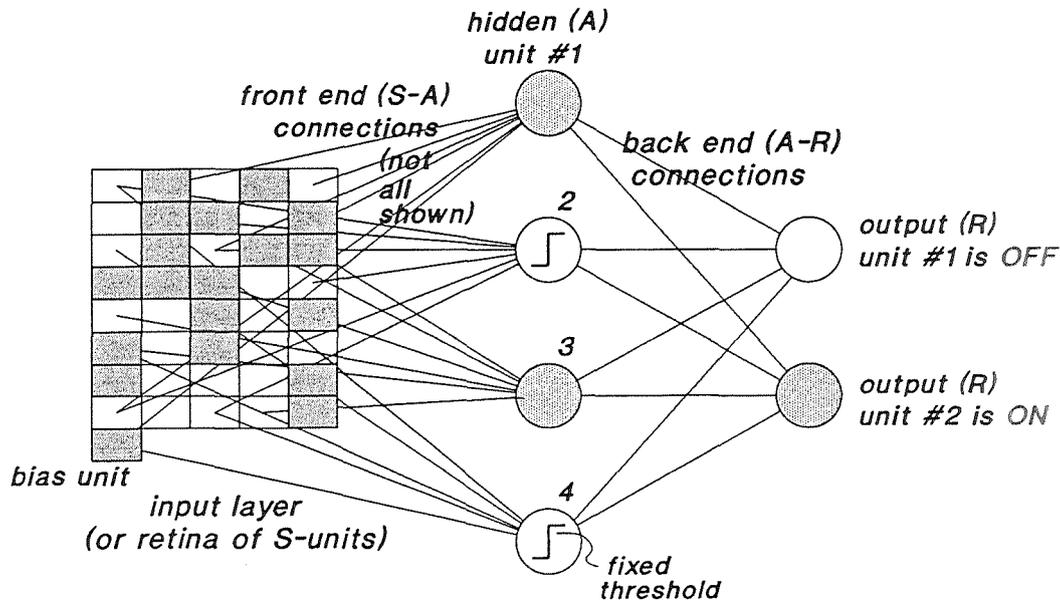


Figure 1. A SIMPLE PERCEPTRON features three layers of neuron-like units with weighted connections feeding excitation (and inhibition) in the forward (l.-to-r.) direction only. Here one of the two output units turns ON to the input pattern.

The revival of perceptron-like models in the 1980s was made possible by a combination of developments, including the widespread perception that AI had reached a plateau, and by the availability of cheaper, faster computers with large amounts of inexpensive RAM (which is needed to store the synaptic weights of large networks). The connectionist models of the 1980s overcame some weaknesses of the Minsky-Papert perceptron. Rumelhart, McClelland and the PDP Research Group (1986), in their first widely-read volume on *Parallel Distributed Processing*, emphasized the importance of having a "hidden layer" of neuron-like units sandwiched between the input and output layers of the network [7]. It was as if Minsky and Papert had done away with the A-units in the perceptron. So these had to be re-invented as "hidden units"! The PDP Group pointed out that these hidden units give three-layer networks the ability--in principle--to solve virtually any pattern classification problem.

But the "powerful new result" that drove the progress of artificial neural networks in the late 1980s was an algorithm called "back-propagation" which permits three-layer networks to learn internal representations of data sets for which no mathematical model can be written down to specify the correct responses to given inputs. Instead, the neural network learns by example in the course of many passes through a training set. In 1986, T. Sejnowski demonstrated NETtalk, the neural network that learned to read aloud in English. The input units in the three-layer network represented sequences of letters from a text. The output units corresponded to the "phonemes" of which spoken English is made. The phonemes were transmitted to a speech synthesizer. NETtalk learned by example to convert letter strings into phonemes. The PDP Group's back-propagation technique was used to modify the weights in a way that resisted and eventually corrected the errors. The speech produced by the network was initially just a meaningless babble. As training progressed around the clock on a mainframe computer, the sounds became more and more intelligible. After the network had learned the training set, it showed the ability to generalize by "reading aloud" the remaining text. This and a legion of other persuasive demos have testified to the power of back-propagation, which has driven the great majority of neural network applications to date. The technology has evolved so far, so fast, that its roots have become almost invisible. According to the prevailing historical view, back-propagation is radically different from the training procedures used with perceptrons [8, 9]. Although the introduction of hidden units gives a feed-forward

network the potential to learn an arbitrary input-to-output mapping, in this view, no technique had existed for training the weights of a network with one or more hidden layers.

It is true that Rosenblatt usually left the weights of the front end (S-to-A) connections at their initial values and applied corrective modifications only to the back end (A-R) weights. In chapter 13 of *Neurodynamics*, however, Rosenblatt addressed the limitations imposed by neglecting to modify the front-weights: "Only one constraint needs to be dropped in order to obtain the most general system of this class: the requirement that the S-to-A connections must have fixed values, only the A-to-R connections being time dependent. In [Chapter 13], variable S-to-A weights will be introduced and the applications of an error-correction procedure will be analyzed. It would seem that a considerable improvement in performance might be obtained if the S-to-A connections could somehow be optimized by a learning process rather than accepting the arbitrary or pre-designed network with which the perceptron starts out. It will be seen that this is indeed the case, provided that certain pitfalls in the design of a reinforcement control procedure are avoided."

With this rationale, Rosenblatt introduced a "back-propagating error correction procedure" consisting of a brief list of rules for assigning errors to hidden (A) units based on their interactions with output (R) units that assume the wrong state in response to the training input. Back-propagation is a "supervised" learning algorithm which obtains its feedback from the output units, computing errors by comparing their observed states to preassigned correct values, propagating errors (and corrections) back towards the front (input) end of the net if a satisfactory solution cannot be found quickly by making corrections at the output end. The actual modification to the weights is formally the same whether an output unit or a hidden unit (or A-unit) is considered. Thus if the error assigned to a unit is positive, the weights of all connections from active units are increased, eventually turning it on. If the error is negative, the weights of connections from active units are decreased. The essential feature of the method is a probabilistic procedure for assigning errors to hidden units.

USING ORIGINAL BACKPROP

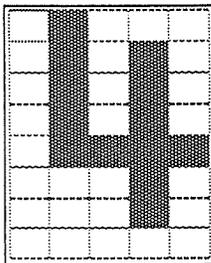
Original Backprop includes four subprograms: (1) *TSET*, a graphical interface for creating the training sets; (2) *BACKPROP*, which generates and trains neural networks; (3) *VIEWNET*, which lets the user analyze and simplify the networks; and (4) *DONET* which recalls and runs the finished product of the neural net design-and-training exercise under other software applications.

TSET presents the user with an 5-by-8 grid of picture elements (or "input units") which can be toggled ON or OFF with a keystroke. In Figure 2, the grid is used to draw 16 patterns representing the hexadecimal symbols zero through F. Pattern number 5, for example, is the symbol "4" which has the binary representation 100. This training set, consisting of four pages of four patterns each, will show an appropriately configured neural net how to map the symbol patterns into binary numbers. Onscreen help is provided for moving around in the pattern set and for naming the individual picture elements when appropriate. In a medical diagnosis problem [1], for example, the picture elements could be placed in one-to-one correspondence with the (40 or fewer) symptoms and named accordingly so that the meaning of the "input unit" is clearly defined as the cursor is moved around the grid in the process of data entry. Training sets are saved as ASCII *.set files. In Figure 2, the file name is *hex.set*; and it is divided into four pages of four patterns each. Version 1.1 limits the size of the training set to 10 pages of binary-valued patterns. Version 1.2 increases the capability to 200 patterns with up to 128 components each and lets the picture elements be represented with 8-bit precision (and 256 colors).

BACKPROP is operated from two menus. The Main Menu presents these Options: (1) get a training set; (2) get a neural net; (3) create a new network; (4) test/train a network; and (5) quit. Option (1) is the obvious starting point. Once a *.set file has been retrieved, it is displayed in a binary string format as shown in Figure 3 for the hex-to-binary conversion problem. The desired mapping is from "input" into "class". The "out" column in the table is all zeros at this point; but one can return to this screen later on (by Option e, below) when training is underway to see how the output units of the network compare to the desired classifications. Selecting Option (3) produces the screen shown in Figure 4. Observe that the numbers of input and output units have defaulted to the numbers indicated by the dimensions of the training set. The number of hidden units has

EDIT/CREATE a Training Set
file name: hex.set

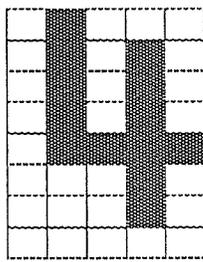
pattern 5



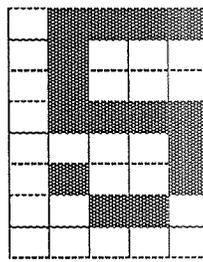
input unit 0:
pattern name: 4

Arrows move cursor; T toggles unit.
Press Insert to record pattern
and advance to next pattern.
Press - (minus) to back up.
F1 = Help

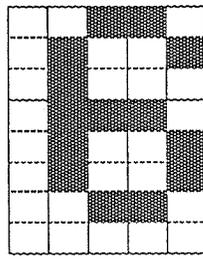
patt. 5 0100



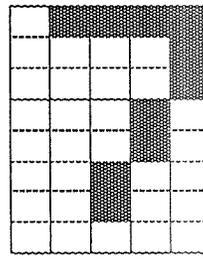
patt. 6 0101



patt. 7 0110

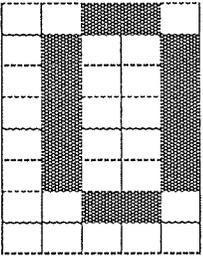


patt. 8 0111

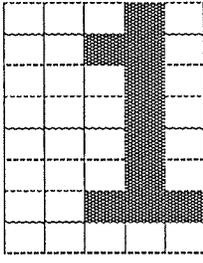


page 2

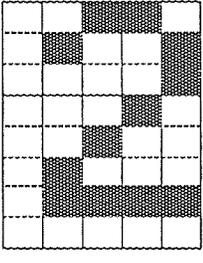
patt. 1 0000



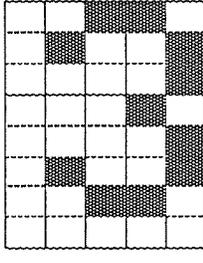
patt. 2 0001



patt. 3 0010

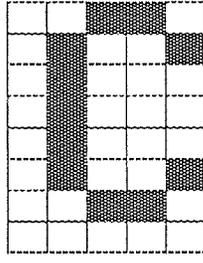


patt. 4 0011

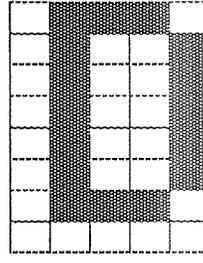


page 1

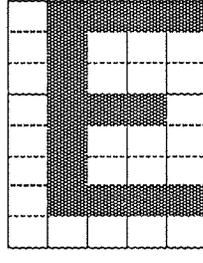
patt. 13 1100



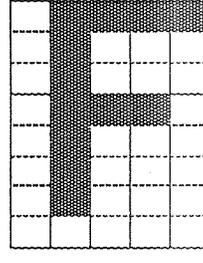
patt. 14 1101



patt. 15 1110



patt. 16 1111



page 4

FIGURE 2. TSET lets Version 1.1 users create input patterns on a 5-by-8 grid and attach as many as five classification bits to each pattern. The product is a training set which is saved as a *.set file.

TRAINING SET

File name: hex.set { 16 patterns }

Page Down to view more patterns. Strike a key to continue.

#	name	input	class	out
1	0	00110010010100101001010010100100110	0000	0000
2	1	00010001100001000010000100001000111	0001	0000
3	2	00110010010000100010001000100001111	0010	0000
4	3	00110010010000100010000010100100110	0011	0000
5	4	01000010100101001010011110001000010	0100	0000
6	5	01111010000100001111000010100100110	0101	0000
7	6	00110010010100001110010010100100110	0110	0000
8	7	01111000010000100010000100010000100	0111	0000
9	8	00110010010100100110010010100100110	1000	0000
10	9	00110010010100100111000010100100110	1001	0000
11	A	00110010010100101111010010100101001	1010	0000
12	B	01110010010100101110010010100101110	1011	0000
13	C	00110010010100001000010000100100110	1100	0000
14	D	01110010010100101001010010100101110	1101	0000
15	E	01111010000100001110010000100001111	1110	0000

FIGURE 3. BACKPROP displays the first 15 elements of the training set (hex.set) as binary strings.

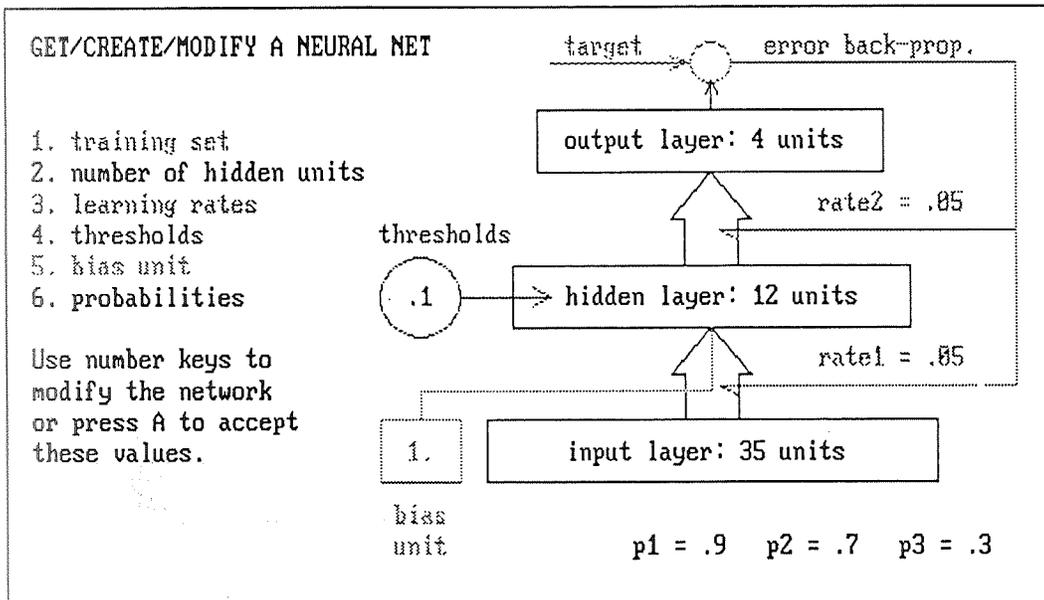


FIGURE 4. BACKPROP creates three-layer networks with the appropriate numbers of input and output units (as required by the training set). This screen serves as a control panel to adjust the number of hidden units and other network parameters.

been set to 12. Thresholds have been established in these hidden units and the input layer bias unit has been turned on.

BACKPROP's Main Menu Option (4) brings up a Training Menu which contains these seven new options: (a) begin training; (b) freeze/unfreeze weights; (c) modify network; (d) save network; (e) review patterns; (f) continue training; and (g) return to Main Menu. Choosing Option (a) now starts the process of learning to associate binary numbers with the symbols in *hex.set*. Figure 5 (top) shows the learning curve which resulted from 278 cycles through the training set of 16 patterns. Two learning curves are actually displayed: A red curve shows the number of incorrectly classified patterns in the current epoch of 50 cycles and a white curve (presently in the upper left corner) shows the average number of errors epoch-by-epoch. It turns out that *hex.set* is a rather difficult assignment. Pressing the ESCape key after cycle number 278, where the error rate has dropped below two-thirds, Option (c) is used to re-access the network parameter control screen of Figure 4. Modifying the "learning rates" (so that *rate1* = .01 and *rate2* = .001), then continuing the training with Option (f), the learning process is rapidly completed as shown in the bottom half of Figure 5. Note that Rosenblatt's stochastic learning algorithm, although it guarantees convergence when a solution exists, does not give the sort of monotonic learning curve that users of PDP back-propagation are accustomed to seeing. The trained network is saved as a *.*net* file after exercising Option (d). Since there are 12 hidden units in the hex-to-binary conversion network, the weights are saved in a file called *hex12.net*.

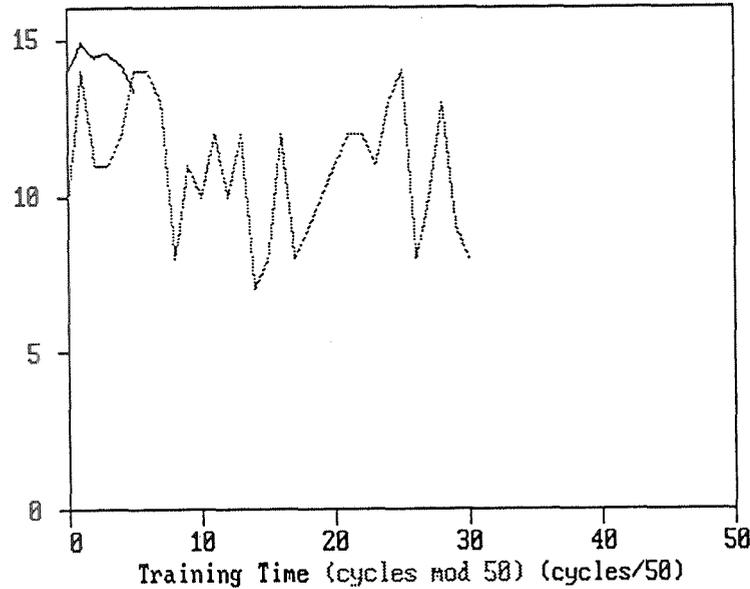
VIEWNET lets the trained network be examined, tested, and "pruned" by the deletion of marginally useful hidden units. *VIEWNET*'s menus control the acquisition of *.*set* and *.*net* files and give the user two "views" of how the network deals with the training set. The detailed view is presented on the "Neural Net Display Screen" (Figure 6) which allocates a small box for each unit and a wide box for each weight. The cursor moves up and down the hidden layer. In Figure 6, the cursor illuminates hidden unit #12; and the weights displayed are all those of the (S-A) connections fanning into this unit from the input layer together with those fanning out to the four output units. On the right side of the screen, the four output units still register "0000" (instead of the desired "0101") because the SPACEbar has not been pressed to propagate the input (pattern #6) forward. A less complicated depiction of the network's performance is obtained by listing the hidden layer activation vectors as columns under the corresponding pattern numbers as in Figure 7.

A recurring question in neural net research concerns the number of hidden units needed to solve the problem presented by the training set. If too few hidden units are employed, training progress may be extremely slow or the solution may actually be unattainable irrespective of any time limit. Use of too many hidden units results in "brittle" solutions and networks that do not generalize well. Some pioneering work of Australian Navy investigators J. Sietsma and R. Dow suggests that the most practical and expedient approach is to first set up and train a network with an abundance of hidden units and then "prune" the trained network by selectively deleting those units which contribute little or nothing to overall performance [10]. *VIEWNET* is the tool that makes it practical to implement such a strategy. From the screen shown in Figure 6, the user can pick a hidden unit (corresponding to a row of the binary array), delete it, and see what effect this has on the correctness of the net's response to each training pattern. Although it requires some work, moving "manually" back-and-forth between *VIEWNET* (to prune one or two units at a time) and *BACKPROP* (to correct the few new errors thus incurred) leads to efficient solutions in much less time than it would take using *BACKPROP* alone. (For example, a network with seven hidden units can be obtained by pruning *hex12.net* in stages; but for *BACKPROP* to solve the problem posed by *hex.set* directly--starting with just seven hidden units--seems to take far more than 250,000 cycles.) A desirable feature which has *not* been included in Version 1.2 (but deferred to later upgrades) is an "autoprune" option which would obviate the need for such "manual" labor.

Version 1.2 improves upon its predecessor by supporting larger training sets and networks. It also includes a new subprogram, *DONET*, to exercise trained networks (retrieved as *.*net* files) and display their responses to given inputs in a dialog box that pops up under other applications (like spreadsheets and word processors). In Version 1.2, *BACKPROP* can be made to find more robust solutions by injecting low-level "noise" into the input patterns in the concluding phases of the training process.

LEARNING CURVE
 training set (file name) = hex.set

Number of Errors (Mean Errors)



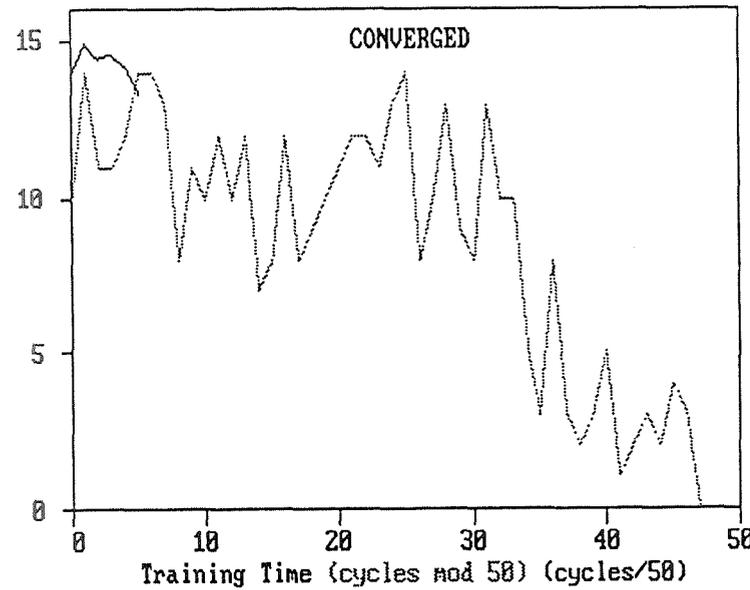
Pattern Number and Status			
correct/incorrect			
1	11	21	31
2	12	22	32
3	13	23	33
4	14	24	34
5	15	25	35
6	16	26	36
7	17	27	37
8	18	28	38
9	19	29	39
10	20	30	40

cycle number 278
 # of errors 8

ESCAPE returns to
 Training Menu

LEARNING CURVE
 training set (file name) = hex.set

Number of Errors (Mean Errors)



Pattern Number and Status			
correct/incorrect			
1	11	21	31
2	12	22	32
3	13	23	33
4	14	24	34
5	15	25	35
6	16	26	36
7	17	27	37
8	18	28	38
9	19	29	39
10	20	30	40

cycle number 294
 # of errors 0

ESCAPE returns to
 Training Menu

FIGURE 5. BACKPROP produces learning curves to show the number of misclassified patterns as a function of the number of cycles through the training set. The slow progress in the first 278 cycles (top) is accelerated by lowering the learning rates from their default values (as described in the text). All 16 patterns in hex.set are learned in 294 cycles (bottom).

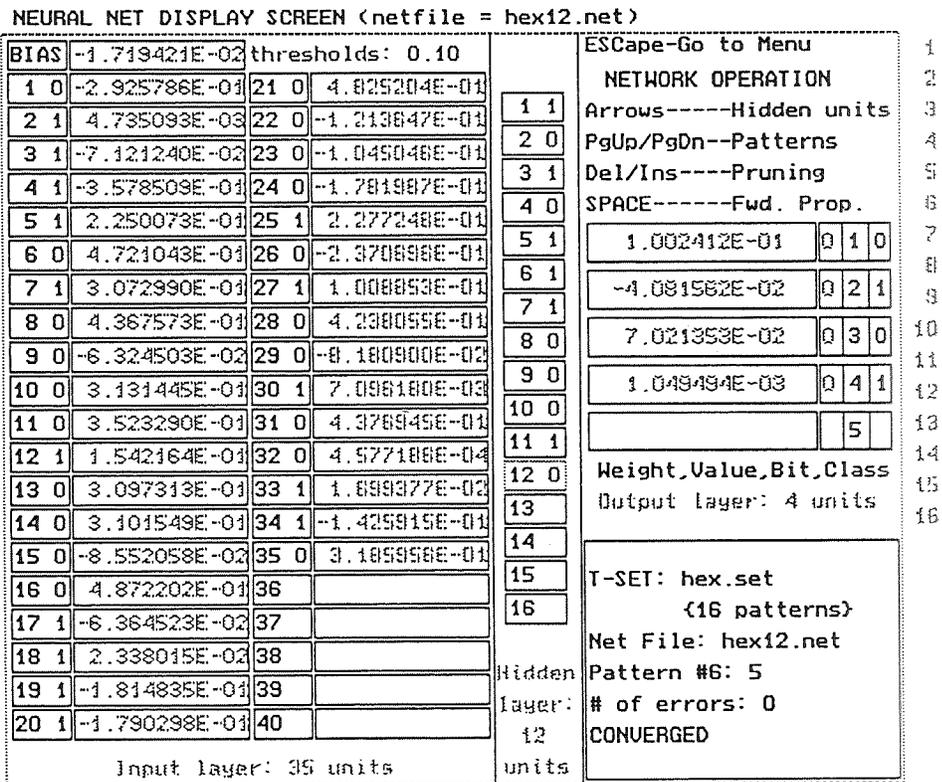


FIGURE 6. VIEWNET's Neural Net Display Screen illuminates the contents of weight file *hex12.net* and shows how the network responds to the elements of the training set.

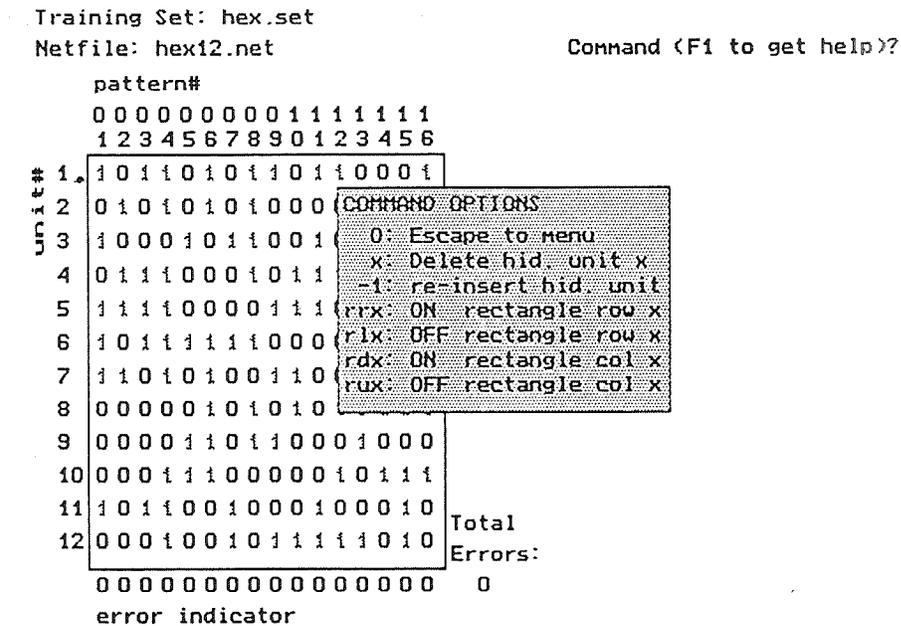


FIGURE 7. VIEWNET gives users the ability to "prune" the network by deleting marginally useful hidden units.

AVAILABILITY

Original Backprop, Version 1.2, will be ready for release in February, 1993. Requests should be sent to the author by regular mail.

ACKNOWLEDGEMENTS

Studies leading to the development of *ORIGINAL BACKPROP* were sponsored by the Office of Naval Research through the Naval Surface Warfare Center's Independent Research Program. Most of the critical components of Version 1.1 were designed and programmed in Borland Turbo-C, during the summer of 1991, by Dovid Lipman, who is presently a student at the Ner Israel Rabbinical College (in Baltimore).

REFERENCES

- [1] Agyei-Mensah, S.O., and Lin, F.C. (1992). Application of neural networks in medical diagnosis: the case of sexually transmitted diseases. *Submitted for publication*.
- [2] Harrison, R.W. (1991). *A neural network for classifying active sonar returns* [Naval Surface Warfare Center, Dahlgren, VA], Tech. Report No. 91-327.
- [3] Coughlin, J.P. (1992). Measures of serial data compressibility by neural network predictors. *Proc. Int'l. Joint Conf. on Neural Nets. [IJCNN Baltimore '92]*, 755-761.
- [4] Rosenblatt, F. (1961). *Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms* [Cornell Aero. Lab., Buffalo, NY], Tech. Report No. VG-1196-G-8].
- [5] Rosenblatt, F. (1962). *Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms* [Spartan Books, Washington, DC].
- [6] Minsky, M.L. and Papert S.A. (1988). *Perceptrons* [Expanded Edn., MIT Press].
- [7] Rumelhart, D.E., McClelland J.L., and the PDP Research Group (1986). *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, Vol. I [MIT Press, Cambridge, MA].
- [8] Widrow, B. and Lehr, M.A. (1990). Thirty years of adaptive neural networks: Perceptron, Madaline, and Backpropagation. *Proceedings of the IEEE* 78(9), 1415-1442
- [9] Denning, P.J. (1992). Neural networks. *American Scientist* 80, 426-429.
- [10] Sietsma, J., and Dow, R.J.F. (1988). Neural net pruning - why and how. *Proc. IEEE Int'l. Conf. on Neural Networks*, 325-333.

CONTROL OF COMPLEX DYNAMIC SYSTEMS BY NEURAL NETWORKS

N 93-25611

350-63

150520

1-10

James C. Spall and John A. Cristion
The Johns Hopkins University
Applied Physics Laboratory
Johns Hopkins Road
Laurel, Maryland 20723-6099 U.S.A.

ABSTRACT

This paper considers the use of neural networks (NN's) in controlling a nonlinear, stochastic system with unknown process equations. The NN is used to model the resulting unknown control law. The approach here is based on using the output error of the system to train the NN controller without the need to construct a separate model (NN or other type) for the unknown process dynamics. To implement such a direct adaptive control approach, it is required that connection weights in the NN be estimated while the system is being controlled. As a result of the feedback of the unknown process dynamics, however, it is not possible to determine the gradient of the loss function for use in standard (back-propagation-type) weight estimation algorithms. Therefore, this paper considers the use of a new stochastic approximation algorithm for this weight estimation, which is based on a "simultaneous perturbation" gradient approximation that only requires the system output error. It is shown that this algorithm can greatly enhance the efficiency over more standard stochastic approximation algorithms based on finite-difference gradient approximations.

1. INTRODUCTION

One of the major problems faced by system designers is finding a means to control and regulate a system when there is uncertainty about the nature of the underlying process. Adaptive control procedures have been developed in a variety of areas for such problems (e.g., robot arm manipulation, materials handling, quality control, etc.), but are typically limited by the need to assume that the forms of the system equations are known (and usually linear) while the parameters may be unknown. In complex physical, socioeconomic, or biological systems, however, the forms of the system equations (typically nonlinear) are often unknown as well as the parameters, making it impossible to determine the control law needed in existing adaptive control procedures. This provides the motivation for considering the use of a neural network (NN) as a controller.

The approach here uses the observed system output error (actual output - target output) to train the NN-based controller without the need to identify or assume a separate model for the system. As we will show, it is not generally possible to train the NN via well-known back-propagation-type algorithms since the required gradient depends on a model for the underlying system. Thus, this paper shows how the simultaneous perturbation stochastic approximation algorithm can be used as a practical weight estimation technique in such a model-free setting. It is shown that this algorithm is much more efficient than more standard finite-difference-based algorithms.

The direct control approach here is based on using a feed-forward NN to approximate the unknown control law. The basis for this approach is the now well-known fact that any measurable function can be

This work was partially supported by the JHU/APL IRAD Program and U.S. Navy Contract N00039-91-C-0001. A more complete version of this paper is available upon request.

approximated to within any degree of accuracy by some single (or multiple) hidden-layer feed-forward NN (e.g., Funahashi [10] or Hornik, Stinchcombe, and White [12]). Our approach will proceed in one of two ways: one method will be based on making almost no assumptions about the nature of the underlying process while the other method will be based on assuming that some information (but still incomplete) is available on the form of the process equations. In the first (basically no structure information) method, the output of the NN will be used to directly approximate the elements of the control vector; in the second (partial structure information), we create a control functional that depends on unknown functions describing the system dynamics and then use a NN to approximate the unknown functions. The second method is reminiscent of the self-tuning regulator approach to adaptive control (e.g., Davis and Vinter [7, pp. 309-312]), except that we are concerned with estimating functions in a control law functional as opposed to estimating parameters in a control law with a fully known structure.

A number of others have considered using NN's for the problem of controlling uncertain nonlinear (usually deterministic) systems (see, e.g., the April 1990 and April 1992 special issues of the IEEE Control Systems Magazine, Narendra and Parthasarathy [22,23], Hunt and Sbarbaro [15], or Iiguni, Sakai, and Tokumaru [16]). Although these methods are useful under certain (fairly restrictive) conditions, they often lack the ability to control systems with minimal prior information. In particular, they require an explicit model (either NN or other parametric type) for the underlying process equations; this model is assumed to be equivalent to the "true" process equations so that it is possible to calculate the gradient needed in back-propagation-type learning algorithms. These techniques (esp. those of Narendra and Parthasarathy) also require off-line identification of the process model before implementation of the adaptive control algorithm. In contrast, our direct control approach uses the NN strictly as a model for use in the control law (no additional NN or parametric model is used for the process directly); the weights in the NN are estimated adaptively based only on the output error of the process (no prior identification is required).¹ As stated in Hoskins, Hwang, and Vagners [13], one of the major advantages of direct control techniques (versus indirect control) is that they are better able to adapt to changes in the underlying system since they are not heavily based on a prior model. Our approach addresses the shortcoming noted in Narendra and Parthasarathy [22, p. 19] that "At present, methods for directly adjusting the control parameters based on the output error (between the plant and [target] outputs) are not available."

Because it is not possible in our framework to obtain the derivatives necessary to implement standard gradient-based search techniques such as back-propagation, we will consider stochastic approximation (SA) algorithms based on approximations to the required gradient. Usually such algorithms are based on standard finite-difference approximations to the gradient (i.e., as in the multivariate Kiefer-Wolfowitz algorithm—see, e.g., Ruppert [26]). These, however, can be very costly in terms of the amount of data required, especially in high-dimensional problems such as estimating a NN weight vector (which easily has dimension of order 10^2 or 10^3). We will, therefore, consider an SA algorithm based on a "simultaneous perturbation" gradient approximation (Spall [29, 30]), which is typically much more efficient than the standard SA algorithms mentioned above in the amount of data required.

The remainder of this paper is organized as follows. Section 2 describes the two related methods for using NN's to control nonlinear systems. This section also describes why it is not possible to determine the gradient of the loss function, in contrast to the approaches of Narendra and others mentioned above where they either assume that the process dynamics is of known structure or introduce an additional NN to model the

¹Chen [4] and Goldenthal and Farrell [11] have also described techniques for NN weight estimation in adaptive control when the gradient is not available, but their techniques have only been developed for particular deterministic model structures and still require considerable information about the process dynamics; in particular they require knowledge of the signs of the terms that appear in the process dynamics. Spall and Cristion [31] includes a more detailed analysis of these techniques.

dynamics. Section 3 discusses the SA approach to weight estimation using a simultaneous perturbation gradient approximation. Section 4 presents a numerical study on a nonlinear system.

2. OVERVIEW OF NEURAL NETWORK APPROACH TO CONTROL

This section describes how the NN will be implemented for the control of uncertain systems. We will describe two methods: one applies when essentially nothing is known about the dynamics of the system and the other applies when partial information on the dynamics is available. The section closes with a discussion of why the well-known "back-propagation" algorithm (or any other algorithm requiring the gradient of the loss function) can not be used for connection weight estimation in this type of general (direct) control problem, which motivates the use of stochastic approximation as discussed in Section 3.

Consider a system output vector at time $k + 1$ given by

$$\mathbf{x}_{k+1} = \Phi_k(\mathbf{f}_k(\mathbf{x}_k, \mathbf{x}_{k-1}, \dots, \mathbf{x}_{k-s}), \mathbf{u}_k, \mathbf{w}_k), \quad s \geq 0, \quad (2.1)$$

where $\Phi_k(\cdot)$ and $\mathbf{f}_k(\cdot)$ are generally unknown, nonlinear functions governing the dynamics of the system, \mathbf{u}_k is the control input applied to affect the system at time $k + 1$, and \mathbf{w}_k is random noise ($\mathbf{f}_k(\cdot)$ may also depend on an arbitrary number of previous controls and/or noise terms, but we omit this generalization for ease of notation). The most important special case of (2.1) is the Markov formulation where $\mathbf{f}_k(\cdot) = \mathbf{f}_k(\mathbf{x}_k)$. Our goal is to choose the sequence of control vectors $\{\mathbf{u}_k\}$ in a manner such that the system output is close to a sequence of target vectors $\{\mathbf{t}_k\}$, where "close" is relative to the magnitude of the noise and the cost associated with the control.

More formally, given the measurements up to time point k we attempt to find the control that minimizes the one-step ahead loss function:

$$L_k \equiv E[(\mathbf{x}_{k+1} - \mathbf{t}_{k+1})^T \mathbf{A}_k (\mathbf{x}_{k+1} - \mathbf{t}_{k+1}) + \mathbf{u}_k^T \mathbf{B}_k \mathbf{u}_k | \mathcal{F}_{k-1}], \quad (2.2)$$

where \mathbf{A}_k , \mathbf{B}_k are positive semi-definite matrices reflecting the relative weight to put on deviations from the target and on the cost associated with larger values of \mathbf{u}_k , and \mathcal{F}_{k-1} is the σ -algebra generated by $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k, \mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{k-1}\}$. An important special case of (2.2) is the minimum variance regulator, where $\mathbf{A}_k = \mathbf{I}$ and $\mathbf{B}_k = \mathbf{0}$.

We will consider two methods to the problem of constructing a controller \mathbf{u}_k in the face of uncertainty about the dynamics of the system, as illustrated in Figs. 2.1a,b for the important special case where $\mathbf{f}_k(\cdot) = \mathbf{f}_k(\mathbf{x}_k)$ (for the more general case as shown in (2.1), the diagrams would be modified in an obvious way). Both of the methods here correspond to direct control approaches as defined (without a solution) in Narendra and Parthasarathy [22] in that NN learning is based directly on the output error, $\mathbf{x}_k - \mathbf{t}_k$; these are in contrast to the indirect control methods of Narendra and Parthasarathy (and others), which are based on the off-line identification of a model of the system based on the error between the system output and model output (not system output and target) for a set of prespecified $\mathbf{0}\mathbf{u}_k$ inputs. In the direct approximation method of Fig. 2.1a, the output of the NN will correspond directly to the elements of the \mathbf{u}_k vector, i.e. the inputs to the NN will be \mathbf{x}_k and \mathbf{t}_{k+1} and the output will be \mathbf{u}_k . This approach is appropriate, e.g., when both $\Phi_k(\cdot)$ and $\mathbf{f}_k(\cdot)$ are unknown functions. In contrast to the direct approximation method of Fig. 2.1a, the NN in the self-tuning method of Fig. 2.1b is used to approximate the unknown dynamics $\mathbf{f}_k(\cdot)$, which is then used in a known functional π_k to obtain \mathbf{u}_k . Since this method requires that π_k (the functional minimizing (2.2)) be known, it requires that the overall relationship between $\mathbf{f}_k, \mathbf{u}_k$ and \mathbf{w}_k , i.e., $\Phi_k(\cdot)$ in (2.1), be known. A very important type of process to which this second method can apply is an affine-nonlinear system as in Chen [4]. When prior information associated with knowledge of $\Phi_k(\cdot)$ is available, the self-tuning method of Fig. 2.1b is often able to yield a superior controller. For both the direct approximation and self-tuning methods, it is required that it be known which arguments appear in $\mathbf{f}_k(\cdot)$, i.e., for the general setting of (2.1) it is required that s be known.

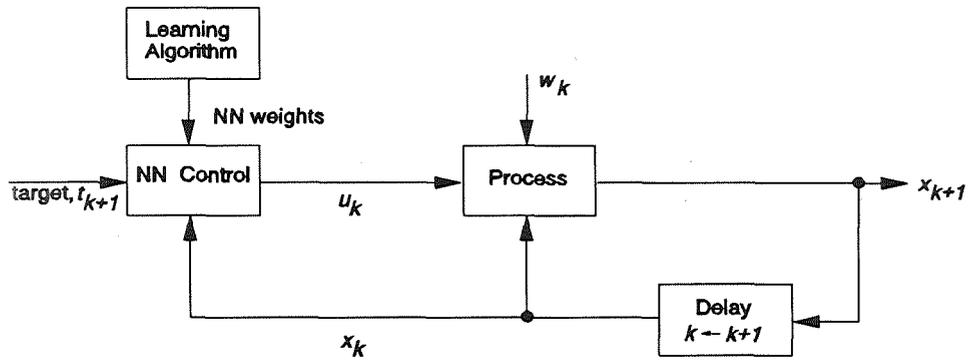


Fig. 2.1a. Control System with NN as Direct Approximator to Optimal u_k when $f_k(\cdot) = f_k(x_k)$

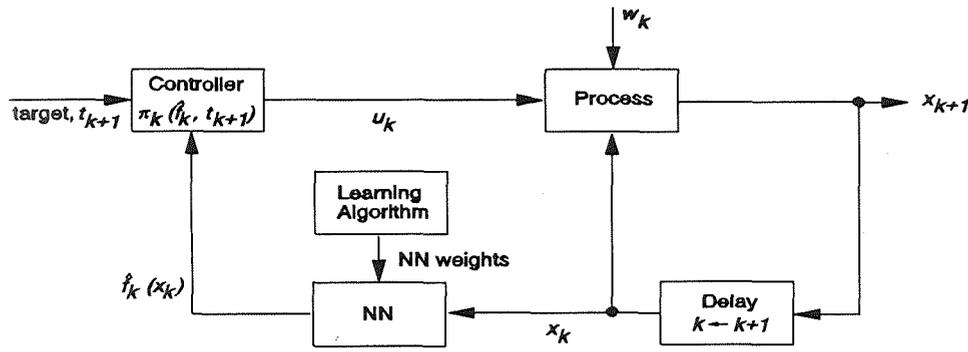


Fig. 2.1b. Self-Tuning Control System with NN as Approximator to $f_k(\cdot)$ when $f_k(\cdot) = f_k(x_k)$

The NN's to be considered here are feed-forward with at least one hidden layer of nodes (neurons) between the input and output nodes. The nodes between (but not within) adjacent layers are all connected and each connection has an associated weight, which is to be estimated from system data. It is this type of NN to which we will restrict our attention, although our method would also apply to other types of NN's (e.g., recurrent). (Since NN's have been discussed in a number of previously published control papers, we will not go into detail here on their development and theory.)

Based on the error criterion in (2.2), we wish to determine the optimal configuration for the NN. Since we assume here that the number of layers and nodes (i.e., network structure) is given, this reduces to a problem of determining the optimal values for the connection weights (determining the NN structure is an important problem in its own right, and has been considered, e.g., in Huang and Huang [14]). Letting $\theta_k \in \mathbb{R}^p$ be the vector of these weights for use in u_k , we are seeking the value of θ_k , say θ_k^* , that minimizes (2.2) given the control as found in Figs 2.1a,b. Thus for each k , we are seeking

$$\left\{ \theta_k^* : \frac{\partial L_k}{\partial \theta_k} = \frac{\partial u_k^T}{\partial \theta_k} \cdot \frac{\partial L_k}{\partial u_k} = 0 \right\}. \quad (2.3)$$

Since $f_k(\cdot)$ (and possibly $\phi_k(\cdot)$) are unknown functions, the term $\partial L_k / \partial u_k$ (and possibly $\partial u_k^T / \partial \theta_k$) in (2.3), which involves the term $\partial \phi_k / \partial u_k$, is not generally computable.² Thus the standard "back-propagation" algorithm (i.e., steepest descent – see, e.g., Narendra and Parthasarathy [23] or White [34]), or any other algorithm involving $\partial L_k / \partial \theta_k$, is not feasible.

To illustrate further why $\partial L_k / \partial \theta_k$ is not available in our direct control setting, consider a simple scalar-deterministic version of system (2.1). Then,

$$\frac{\partial L_k}{\partial \theta_k} = \frac{\partial (x_{k+1} - t_{k+1})^2}{\partial \theta_k} = 2(x_{k+1} - t_{k+1}) \frac{\partial \phi_k}{\partial u_k} \frac{\partial u_k}{\partial \theta_k}. \quad (2.4)$$

When neither $\phi_k(\cdot)$ nor $f_k(\cdot)$ is known (as in the direct approximation method of Fig. 2.1a), then neither of the derivatives on the right-hand side of (2.4) will be known. When $f_k(\cdot)$ is unknown and $\phi_k(\cdot)$ and $u_k = \pi_k(\cdot)$ are known (as in the self-tuning method of Fig. 2.1b), then $\partial u_k / \partial \theta_k$ will be known but $\partial \phi_k / \partial u_k$ will, in general, still be unknown since it will depend on $f_k(\cdot)$.³ Thus we see that in either of the direct control settings in Fig. 2.1a,b, $\partial L_k / \partial \theta_k$ is not generally available. The same principles apply in the more general multivariate stochastic version of model (2.1).

Because back-propagation-type algorithms are not generally feasible in the direct control setting here, we consider a stochastic approximation (SA) algorithm of the form

$$\hat{\theta}_k = \hat{\theta}_{k-1} - a_k (\text{gradient approx.})_k \quad (2.5)$$

to estimate $\{\theta_k^*\}$, where $\hat{\theta}_k$ denotes the estimate at the given iteration, $\{a_k\}$ is a scalar gain sequence satisfying certain regularity conditions, and the gradient approximation is such that it does not require knowledge of $f_k(\cdot)$ (and $\phi_k(\cdot)$, if appropriate). The next section is devoted to describing in more detail the SA approach to this problem.

3. WEIGHT ESTIMATION BY SIMULTANEOUS PERTURBATION STOCHASTIC APPROXIMATION

Recall that we are seeking the NN weight vector at each time point that minimizes (2.2), i.e., we are seeking the minimizing θ_k, θ_k^* , such that

$$g_k(\theta_k^*) \equiv \left. \frac{\partial L_k}{\partial \theta_k} \right|_{\theta_k^*} = 0,$$

where θ_k is for use in the control u_k . Recall also that since back-propagation (or other derivative-based) algorithms are not applicable, we will consider an SA-based approach. This subsection describes the simultaneous perturbation SA (SPSA) approach to this problem and mentions how this approach contrasts with the more standard finite-difference SA (FDSA) approach of Kiefer-Wolfowitz. Spall [30] gives a detailed analysis of the SPSA approach to optimization. It is shown that the SPSA algorithm can achieve the same level of

² This contrasts with the "open loop" identification problems in, e.g., Narendra and Parthasarathy [22, Sect. 5], where in estimating the connection weights no unknown functions appear in the gradient. This also contrasts with implementations of so-called indirect feedback controllers (e.g., Narendra and Parthasarathy [22, Sect. 6]) where a NN is used to model the unknown system dynamics and the identification and adaptive control is performed as if the NN model was identical in structure to the true system dynamics.

³ One special case where $\partial L_k / \partial \theta_k$ can be computed is in the self-tuning setting of Fig. 2.1b where $u_k(\cdot)$ is known to enter $\phi_k(\cdot)$ additively (since $\partial \phi_k / \partial u_k$ then does not depend on $f_k(\cdot)$). Of course, in the more general setting of direct approximation control (Fig. 2.1a) $\partial L_k / \partial \theta_k$ would still be unavailable.

asymptotic accuracy as FDSA with only $1/\rho$ the number of system measurements. This is of particular interest in neural network problems since ρ can easily be on the order of 10^2 or 10^3 .

In line with (2.5), the SPSA algorithm has the form

$$\hat{\theta}_k = \hat{\theta}_{k-1} - a_k \hat{g}_k(\hat{\theta}_{k-1}) \quad (3.1a)$$

where $\hat{g}_k(\hat{\theta}_{k-1})$ is the simultaneous perturbation approximation to $g_k(\hat{\theta}_{k-1})$. In particular the ℓ^{th} component of $\hat{g}_k(\hat{\theta}_{k-1})$, $\ell = 1, 2, \dots, \rho$, is given by

$$\hat{g}_{\ell}(\hat{\theta}_{k-1}) = \frac{\hat{J}_k^{(+)} - \hat{J}_k^{(-)}}{2c_k \Delta_{k\ell}}, \quad (3.1b)$$

where

- $\hat{J}_k^{(\pm)} = (x_{k+1}^{(\pm)} - t_{k+1})^T A_k (x_{k+1}^{(\pm)} - t_{k+1}) + u_k^{(\pm)T} B_k u_k^{(\pm)}$,
- $u_k^{(\pm)} = u_k(x_k, \dots, x_{k-s}, t_{k+1}, \hat{\theta}_{k-1} \pm c_k \Delta_k)$, i.e., a control based on a NN with weight vector $\theta_k = \hat{\theta}_{k-1} + c_k \Delta_k$ or $\theta_k = \hat{\theta}_{k-1} - c_k \Delta_k$,
- $x_{k+1}^{(\pm)}$ is system output based on $u_k^{(\pm)}$, $\Delta_k = (\Delta_{k1}, \Delta_{k2}, \dots, \Delta_{k\rho})^T$, with the $\{\Delta_{ki}\}$ independent, symmetrically distributed (about 0) random variables $\forall k, i$, identically distributed at each k , with $E(\Delta_{ki}^2)$ uniformly bounded $\forall k, i$,
- $\{c_k\}$ is a sequence of positive numbers satisfying certain regularity conditions.

The key fact to observe is that at any iteration only two measurements are needed (i.e., the numerators are the same for all ρ components). This is in contrast to the standard FDSA approach where 2ρ measurements are needed to construct the approximation to $g_k(\cdot)$ (i.e., for the ℓ^{th} component of the gradient approximation, the quantity Δ_k is replaced by a vector with a positive constant in the ℓ^{th} place and zeroes elsewhere; see, e.g., Ruppert [26]). A variation on the form in (3.1b) is to average several gradient approximations, with each vector in the average being based on a new (independent) value of Δ_k and a corresponding new pair of measurements; this will often enhance the performance of the algorithm as illustrated in Section 4. A further variation on (3.1b) is to smooth across time by a weighted average of the previous and current gradient estimates (analogous to the "momentum" approach in back-propagation); such smoothing can often improve the performance of the algorithm (see Spall and Cristion [32] for a thorough discussion of smoothing in SPSA-based direct adaptive control).

The complete version of this paper gives a much fuller account of the theory behind SPSA together with some of the practical issues associated with its implementation in adaptive control.

4. EMPIRICAL STUDY

4.1 Preliminaries

This section presents the results of our study on a stochastic generalization of a model in Narendra and Parthasarathy [22] (N & P hereafter in this section). We will compare the SPSA and FDSA weight estimation algorithms.

The study here is based on $A_k = I$ and $B_k = \mathbf{0}$ in the loss function (2.2) (i.e., a minimum variance regulator). The performance of the various techniques will be evaluated by comparing the root-mean-square (RMS) tracking error as normalized by the dimension of x_k , i.e. RMS at time k is

$[(x_k - t_k)^T(x_k - t_k)/\dim(x_k)]^{1/2}$. The (feedforward) NN's considered here have an input layer, two hidden layers, and an output layer, as in N & P and Chen [4]. The hidden layer nodes are hyperbolic tangent functions (i.e., $(e^x - e^{-x})/(e^x + e^{-x})$ for input x) while the output nodes are linear functions (simply x). Each node takes as an input (x) the weighted sum of outputs of all nodes in the previous layer plus a bias weight not connected to the rest of the network (hence an $N_{4,20,10,2}$ network, in the notation of N & P, has $100 + 210 + 22 = 332$ weights to be estimated). For the weight estimation, we will consider different forms of the SPSA algorithm, denoted SPSA- q , where q denotes the number of individual gradient approximations of the form (3.1b) that are averaged to form $\hat{g}_k(\cdot)$ (hence an SPSA- q algorithm uses $2q$ measurements to form $\hat{g}_k(\cdot)$). For the SPSA algorithms we take the perturbations Δ_{k_i} to be Bernoulli ± 1 distributed, which satisfies the relevant regularity conditions of Section 3.

4.2 Results of Numerical Study

The model we consider is a generalization of the two-dimensional model with additive control given in N & P to include additive (independent) noise, i.e.,

$$x_{k+1} = f(x_k) + u_k + w_k, \quad w_k \sim N(0, \sigma^2 I), \quad (4.1)$$

where, as in eqn. (18) of N & P, the data are generated according to

$$f(x_k) = \frac{1}{1 + x_{k2}^2} \begin{bmatrix} x_{k1} \\ x_{k1}x_{k2} \end{bmatrix}$$

with x_{k_i} the i^{th} ($i = 1, 2$) component of x_k . Analogous to N & P we take the two-dimensional target sequence to be generated by the deterministic difference equation

$$t_{k+1} = \begin{pmatrix} .6 & .2 \\ .1 & -.8 \end{pmatrix} t_k + \begin{pmatrix} \sin(2\pi k/25) \\ \cos(2\pi k/25) \end{pmatrix}, \quad t_0 = 0.$$

Because (4.1) is an additive control model, we will only consider the DA method in this study (see footnote 3). To implement DA the analyst is assumed to know that $s=0$ (i.e., that this is a Markov-type model) and that $\dim u_k = 2$. As with N & P we used NN's with two hidden layers, one of 20 nodes and one of 10 nodes (so an $N_{4,20,10,2}$ network was used for the controller). The indicated RMS errors throughout this study are normalized for the two-dimensional setting as discussed in Subsection 4.1; therefore, since $\text{cov}(w_k) = \sigma^2 I$, we know that long-run RMS can at best equal σ .

Fig. 4.1 presents the main results for our study of the model in (4.1). The RMS curves in the figures are based on the sample mean of four independent runs with different initial weights $\hat{\theta}_0$, where the elements of $\hat{\theta}_0$ were generated randomly from a uniform $(-1, .1)$ distribution. To effect a fair comparison of the algorithms the same four sets of initial weight vectors were used for the three different curves. To further smooth the resulting error curves and to show typical performance (not just case-dependent variation), we applied the MATLAB low-pass interpolating function INTERP to the error values based on the average of four runs. The curves shown in the figures are based on this combination of across-realization averaging and across-iteration interpolation. Each of the curves was generated by using SA gains of the form $a_k = A/k^{.7501}$ and $c_k = C/k^{.25}$ with $A, C > 0$ (the exponents were chosen to satisfy standard SA conditions and to afford a_k the effectively slowest rate of decay consistent with these conditions; a slow decay rate tended to accelerate the rate of decrease in RMS error). For each curve we attempted to tune A and C to maximize the rate of convergence of $\hat{\theta}_k$ (as would typically be done in practice); the values satisfied $.037 \leq A \leq .12$ and $.20 \leq C \leq .25$. The value x_0 was set to $(1.5, 1.5)^T$ for all studies, so the initial RMS error is 1.5.

Fig. 4.1 shows that both the SPSA and FDSA algorithms yield controllers with decreasing RMS tracking error over time. The RMS error curves for both algorithms show the characteristic shape of first order (steepest-descent-type) algorithms in that there is a sharp initial decline followed by slow decline. We see that the long-run performance of SPSA-4 is slightly better than that of FDSA, with SPSA-4 and FDSA achieving terminal RMS errors of .32 and .33 respectively (vs. the theoretical limit of .25); for the SPSA-1 algorithm the terminal error was .47. The critical observation to make here is that the SPSA algorithms achieved their performance with a large savings in data: each iteration of SPSA-1 and SPSA-4 used only three measurements and nine measurements, respectively, while each iteration of FDSA used 665 measurements (these measurement counts include the one operational measurement for each iteration in addition to the measurements generated for purposes of constructing the gradient approximation). Hence Fig. 4.1 illustrates that SPSA-4 yields a slightly lower level of long-run tracking error than the standard FDSA algorithm with a 74-fold savings in system measurements. The data savings seen in Fig. 4.1 is typical of that for a number of other studies involving SPSA and FDSA that we have conducted on model (4.1) as well as on other nonlinear models; in fact even greater data savings are typical with more complex NN's (as might be needed in higher-dimensional systems or in systems where u_k is not simply additive).

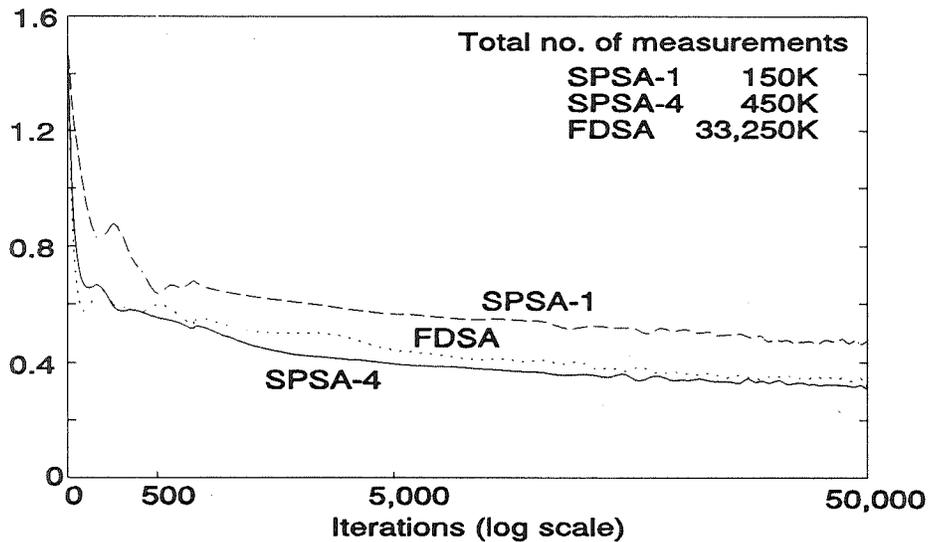


Fig. 4.1 RMS Error for DA Controller with SPSA and FDSA Algorithms in Additive Control Model with $\sigma = .25$

Our other numerical study for (4.2) illustrates the relative performance of SPSA and FDSA in a deterministic ($\sigma=0$) system (to complement the stochastic comparison in the study for (4.1)). As with (4.1), we used the DA control method. The mean and terminal RMS errors for SPSA-1 were, respectively, .12 and .087 versus .14 and .103 for FDSA. Thus SPSA outperforms FDSA with less than 1/220 the number of system measurements.

REFERENCES

(The list below includes all items cited in the full version of this paper, which is available from the authors upon request.)

- [1] Bayard, D.S. [1991], "A Forward Method for Optimal Stochastic Nonlinear and Adaptive Control," IEEE Trans. Auto. Control, vol. 36, pp. 1046-1053.
- [2] Benveniste, A. and Ruget, G. [1982], "A Measure of the Tracking Capability of Recursive Stochastic Algorithms with Constant Gains," IEEE Trans. Auto. Control, vol. AC-27, pp. 639-649.
- [3] Boguslavskij, I.A [1988], Filtering and Control, Optimization Software Pubs. Div., New York.
- [4] Chen, F.C. [1990], "Back-Propagation Neural Networks for Nonlinear Self-Tuning Adaptive Control," IEEE Control Syst. Mag., April, pp. 44-48.
- [5] Chin, D.C. [1992], "A More Efficient Global Optimization Algorithm Based on Styblinski and Tang (1990)," Neural Nets., submitted.
- [6] Chung, K.L. [1974], A Course in Probability Theory, Academic, New York.
- [7] Davis, M.H.A and Vinter, R.B. [1985], Stochastic Modeling and Control, Chapman and Hall, New York.
- [8] Evans, S.N. and Weber, N.C. [1986], "On the Almost Sure Convergence of a General Stochastic Approximation Procedure," Bull. Austral. Math. Soc., vol. 34, pp. 335-342.
- [9] Fabian, V. [1971], "Stochastic Approximation," in Optimizing Methods in Statistics (J.J. Rustagi, ed.), Academic, New York, pp. 439-470.
- [10] Funahashi, K.I. [1989], "On the Approximate Realization of Continuous Mappings by Neural Networks," Neural Nets., vol. 2, pp. 183-192.
- [11] Goldenthal, W. and Farrell, J. [1990], "Application of Neural Networks to Automatic Control," Proc. AIAA Guidance, Navigation and Control Conf., Part 2, pp. 1108-1112.
- [12] Hornik, K., Stinchcombe, M. and White, H. [1989], "Multilayer Feedforward Networks are Universal Approximators," Neural Nets., vol. 2, pp. 359-366.
- [13] Hoskins, D.A., Hwang, J.N., and Vagners, J. [1992], "Iterative Inversion of Neural Networks and its Applications to Adaptive Control," IEEE Trans. Neural Nets., vol. 3, pp. 292-301.
- [14] Huang, S.-C. and Huang, Y.-F. [1991], "Bounds on the Number of Hidden Neurons in Multilayer Perceptrons," IEEE Trans. Neural Nets., vol. 2, pp. 47-55.
- [15] Hunt, K.J. and Sbarbaro, P. [1991], "Neural Networks for Nonlinear Model Control," IEEE Proc.-D, vol. 138, pp. 431-438.
- [16] Iiguni, Y, Sakai, H. and Tokumaru, H. [1991], "A Nonlinear Regulator Design in the Presence of System Uncertainties Using Multilayered Neural Networks," IEEE Trans. Neural Nets., vol. 2, pp. 410-417.
- [17] Isidori, A. and Byrnes, C.I. [1990], "Output Regulation of Nonlinear Systems," IEEE Trans. Auto. Control, vol. 35, pp. 131-140.

- [18] Kuan, C.-M. and Hornik, K. [1991], "Convergence of Learning Algorithms with Constant Learning Rates," IEEE Trans. Neural Nets., vol. 2, pp. 484-489.
- [19] Kushner, H. J. and Huang, H. [1981], "Asymptotic Properties of Stochastic Approximations with Constant Coefficients," SIAM J. Control Optimiz., vol. 19, pp. 87-105.
- [20] Macchi, O. and Eweda, E. [1983], "Second Order Convergence Analysis of Stochastic Adaptive Linear Filtering," IEEE Trans. Auto. Control, vol. 28, pp. 76-85.
- [21] Moden, P.E. and Soderstrom, T. [1982], "Stationary Performance of Linear Stochastic Systems Under Single Step Optimal Control," IEEE Trans. Auto. Control, vol. AC-27, pp. 214-216.
- [22] Narendra, K.S. and Parthasarathy, K. [1990], "Identification and Control of Dynamical Systems Using Neural Networks," IEEE Trans. Neural Nets. vol. 1, pp. 4-26.
- [23] Narendra, K.S. and Parthasarathy, K. [1991], "Gradient Methods for the Optimization of Dynamic Systems Containing Neural Networks," IEEE Trans. Neural Nets., vol. 2, pp. 252-262.
- [24] Nijmeijer, H. and van der Schaft, A. J. [1990], Nonlinear Dynamical Control Systems, Springer-Verlag, New York.
- [25] Psaltis, D., Athanasios, S., and Yamamura, A.A. [1988], "A Multilayered Neural Network Controller," IEEE Control Syst. Mag., vol. 8, April, pp. 17-21.
- [26] Ruppert, D. [1983], "Kiefer-Wolfowitz Procedure," Encyclopedia of Statistical Science, vol. 4 (S. Kotz and N.L. Johnson, eds.), pp. 379-381, Wiley, New York.
- [27] Ruppert, D. [1985], "A Newton-Raphson Version of Multivariate Robbins-Monro Procedure," Ann. Stat., vol. 13, pp. 236-245.
- [28] Saridis, G.N. [1977], Self-Organizing Control of Stochastic Systems, Marcel Dekker, New York.
- [29] Spall, J.C. [1988], "A Stochastic Approximation Algorithm for Large-Dimensional Systems in the Kiefer-Wolfowitz Setting," Proc. IEEE Conf. Dec. Control, pp. 1544-1548.
- [30] Spall, J.C. [1992], "Multivariate Stochastic Approximation Using a Simultaneous Perturbation Gradient Approximation," IEEE Trans. Auto. Control, vol. 37, pp. 332-341.
- [31] Spall, J.C. and Cristion, J.A. [1991], "Neural Networks for Control of Uncertain Systems," Proc. Test Technology Symp. IV (sponsored by U.S. Army Test and Evaluation Command), pp. 575-588.
- [32] Spall, J.C. and Cristion, J.A. [1992], "Nonlinear Adaptive Control Using Neural Networks: Estimation Based on a Smoothed Form of Simultaneous Perturbation Gradient Approximation," Statistica Sinica, submitted.

ADAPTIVE PROCESS CONTROL WITH FUZZY LOGIC AND GENETIC ALGORITHMS

C. L. Karr
Mechanical Engineer
U.S. Bureau of Mines, Tuscaloosa Research Center
P.O. Box L, University of Alabama Campus
Tuscaloosa, AL 35486-9777

551-63
150521
p. 8

ABSTRACT

Researchers at the U.S. Bureau of Mines have developed adaptive process control systems in which genetic algorithms (GAs) are used to augment fuzzy logic controllers (FLCs). GAs are search algorithms that rapidly locate near-optimum solutions to a wide spectrum of problems by modeling the search procedures of natural genetics. FLCs are rule based systems that efficiently manipulate a problem environment by modeling the "rule-of-thumb" strategy used in human decision-making. Together, GAs and FLCs possess the capabilities necessary to produce powerful, efficient, and robust adaptive control systems. To perform efficiently, such control systems require a *control element* to manipulate the problem environment, an *analysis element* to recognize changes in the problem environment, and a *learning element* to adjust to the changes in the problem environment. Details of an overall adaptive control system are discussed. A specific laboratory acid-base pH system is used to demonstrate the ideas presented.

INTRODUCTION

The need for efficient process control has never been more important than it is today because of economic stresses forced on industry by processes of increased complexity and by intense competition in a world market. No industry is immune to the cost savings necessary to remain competitive; even traditional industries such as mineral processing [1], chemical engineering [2], and wastewater treatment [3] have been forced to implement cost-cutting measures. Cost-cutting generally requires the implementation of emerging techniques that are often more complex than established procedures. The new processes that result are often characterized by rapidly changing process dynamics. Such systems prove difficult to control with conventional strategies, because these strategies lack an effective means of adapting to change. Furthermore, the mathematical tools employed for process control can be unduly complex even for simple systems.

In order to accommodate changing process dynamics yet avoid sluggish response times, adaptive control systems must alter their control strategies according to the current state of the process. Modern technology in the form of high-speed computers and artificial intelligence (AI) has opened the door for the development of control systems that adopt the approach to adaptive control used by humans, and perform more efficiently and with more flexibility than conventional control systems. Two powerful tools for adaptive control that have emerged from the field of AI are fuzzy logic [4] and genetic algorithms (GAs) [5].

The U.S. Bureau of Mines has developed an approach to the design of adaptive control systems, based on GAs and FLCs, that is effective in problem environments with rapidly changing dynamics. Additionally, the resulting controllers include a mechanism for handling inadequate feedback about the state or condition of the problem environment. Such controllers are more suitable than past control systems for recognizing, quantifying, and adapting to changes in the problem environment.

The adaptive control systems developed at the Bureau of Mines consist of a *control element* to manipulate the problem environment, an *analysis element* to recognize changes in the problem environment, and a *learning element* to adjust to the changes in the problem environment. Each component employs a GA, a FLC, or both, and each is described in this paper. A particular problem environment, a laboratory acid-base pH system, serves as a forum for presenting the details of a Bureau-developed, adaptive controller. Preliminary results are presented to demonstrate the effectiveness of a GA-based FLC for each of the three individual elements. Details of the system will appear in a report by Karr and Gentry [6].

PROBLEM ENVIRONMENT

In this section, a pH system is introduced to serve as a forum for presenting the details of a stand-alone, comprehensive, adaptive controller developed at the U.S. Bureau of Mines; emphasis is on the method not the application. The goal of the control system is to drive the pH to a setpoint. This is a non-trivial task since the pH system contains both nonlinearities and changing process dynamics. The nonlinearities occur because the output of pH sensors is proportional to the logarithm of hydrogen ion concentration. The source of the changing process dynamics will be described shortly.

A schematic of the pH system under consideration is shown in Fig. 1. The system consists of a beaker and five, valved input streams. The beaker initially contains a given volume of a solution having some known pH. The five, valved input streams into the beaker are divided into the two *control input streams* and the three *external input streams*. Only the valves associated with the two control input streams can be adjusted by the controller. Additionally, as a constraint on the problem, these valves can only be adjusted a limited amount (0.5 mL/s, which is 20 pct of the maximum flow rate of 2.5 mL/s) to restrict pressure transients in the associated pumping systems.

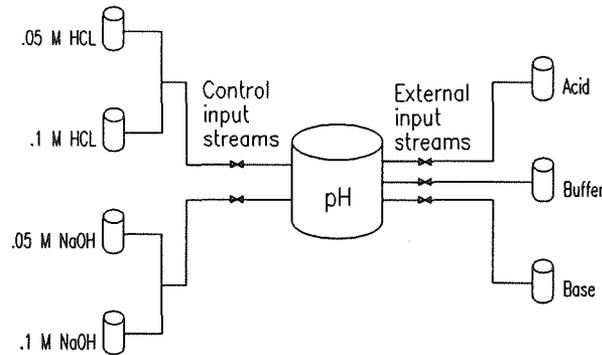


Fig. 1. Basic structure of the pH system.

The goal of the control problem is to drive the system pH to the desired setpoint in the shortest time possible by adjusting the valves on the two control input streams. Achieving this goal is made considerably more difficult by incorporating the potential for changing the process dynamics. These changing process dynamics come from three random changes that can be made to the pH system. First, the concentrations of the acid and base of the two control input streams can be changed randomly to be either 0.1 M HCl or 0.05 M HCl and 0.1 M NaOH or 0.05 M NaOH. Second, the valves on the external input streams can be randomly altered. This allows for the external addition of acid (0.05 M HCl), base (0.05 M CH_3COONa), and buffer (a combination of 0.1 M CH_3COOH and 0.1 M CH_3COONa) to the pH system. Note that the addition of a buffer is analogous to adding inertia to a mechanical system. Third, random changes are made to the setpoint to which the system pH is to be driven. These three random alterations in the system parameters dramatically alter the way in which the problem environment reacts to adjustments made by the controller to the valves on the control input streams. Furthermore, the controller receives no feedback concerning these random changes.

The pH system was designed on a small scale so that experiments could be performed in limited laboratory space. Titrations were performed in a 1,000-mL beaker using a magnetic bar to stir the solution. Peristaltic pumps were used for the five input streams. An industrial pH electrode and transmitter sent signals through an analog-to-digital board to a 33-MHz 386 personal computer which implemented the control system.

STRUCTURE OF THE ADAPTIVE CONTROLLER

Figure 2 shows a schematic of the Bureau's adaptive control system. The heart of this control system is the loop consisting of the control element and the problem environment. The control element receives information from sensors in the problem environment concerning the status of the *condition variables*, i.e., pH and ΔpH . It then computes a desirable state for a set of *action variables*, i.e., flow rate of acid (Q_{ACID}) and flow rate of base (Q_{BASE}). These changes in the action variables force the problem environment toward the setpoint. This is the basic approach adopted for the design of virtually any closed loop control system, and in and of itself includes no mechanism for adaptive control.

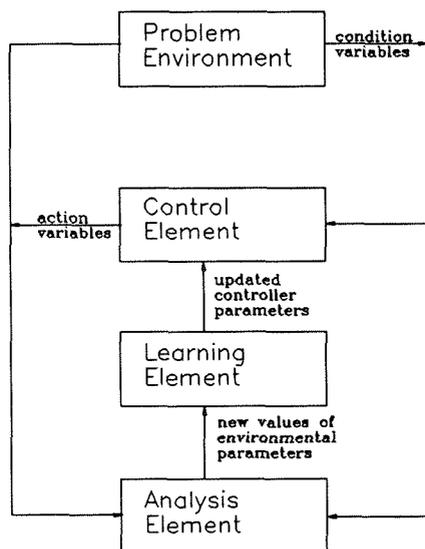


Fig. 2. Structure of the adaptive control system.

The adaptive capabilities of the system shown in Fig. 2 are due to the analysis and learning elements. In general, the analysis element must recognize when a change in the problem environment has occurred. A "change," as it is used here, consists of any of the three random alterations to a parameter possible in the problem environment. (Of importance is the fact that all of these changes affect the response of the problem environment, otherwise it has no effect on the way in which the control element must act to efficiently manipulate the problem environment.) The analysis element uses information concerning the condition and action variables over some finite time period to recognize changes in the environment and to compute the new performance characteristics associated with these changes.

The new environment (the problem environment with the altered parameters) can pose many difficulties for the control element, because the control element is no longer manipulating the environment for which it was designed. Therefore, the algorithm that drives the control element must be altered. As shown in the schematic of Fig. 2, this task is accomplished by the learning element. The most efficient approach for the learning element to use to alter the control element is to utilize information concerning the past performance of the control system. The strategy used by the control, analysis, and learning elements of the stand-alone, comprehensive adaptive controller being developed by the U.S. Bureau of Mines is provided in the following sections.

Control Element

The control element receives feedback from the pH system, and based on the current state of pH and ΔpH , must prescribe appropriate values of Q_{ACID} and Q_{BASE} . Any of a number of closed-loop controllers could be used for this

element. However, because of the flexibility needed in the control system as a whole, a FLC is employed. Like conventional rule-based systems (expert systems), FLCs use a set of production rules which are of the form:

IF {*condition*} THEN {*action*}

to arrive at appropriate control actions. The left-hand-side of the rules (the *condition* side) consists of combinations of the controlled variables (pH and ΔpH); the right-hand-side of the rules (the *action* side) consists of combinations of the manipulated variables (Q_{ACID} and Q_{BASE}). Unlike conventional expert systems, FLCs use rules that utilize fuzzy terms like those appearing in human rules-of-thumb. For example, a valid rule for a FLC used to manipulate the pH system is:

IF {ph is **VERY ACIDIC** and ΔpH is **SMALL**} THEN { Q_{BASE} is **LARGE** and Q_{ACID} is **ZERO**}.

This rule says that if the solution is very acidic and is not changing rapidly, the flow rate of the base should be made to be large and the flow rate of the acid should be made to be zero.

The fuzzy terms are subjective; they mean different things to different "experts," and can mean different things in varying situations. Fuzzy terms are assigned concrete meaning via fuzzy membership functions [4]. The membership functions used in the control element to describe pH appear in Fig. 3. (As will be seen shortly, the learning element is capable of changing these membership functions in response to changes in the problem environment.) These membership functions are used in conjunction with the rule set to prescribe single, crisp values of the action variables (Q_{ACID} and Q_{BASE}). Unlike conventional expert systems, FLCs allow for the enactment of more than one rule at any given time. The single crisp action is computed using a weighted averaging technique that incorporates both a *min-max* operator and the *center-of-area* method [7]. The following fuzzy terms were used, and therefore "defined" with membership functions, to describe the significant variables in the pH system:

- pH Very Acidic (VA), Acidic (A), Mildly Acidic (MA), Neutral (N), Mildly Basic (MB), Basic (B), and Very Basic (VB);
- ΔpH Small (S) and Large (L);
- Q_{ACID} Zero (Z), Very Small (VS);
- Q_{BASE} Small (S), Medium (M), and Large (L).

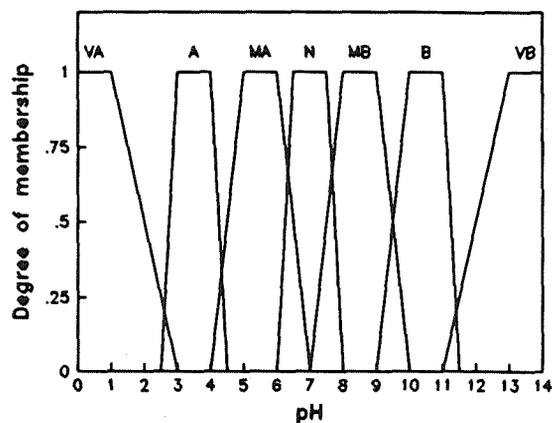


Fig. 3. pH membership functions.

Although the pH system is quite complex, it is basically a titration system. An effective FLC for performing titrations can be written that contains only 14 rules. The 14 rules are necessary because there are seven fuzzy terms

describing pH and two fuzzy terms describing ΔpH ($7 \times 2 = 14$ rules to describe all possible combinations that could exist in the pH system as described by the fuzzy terms represented by the membership functions selected). Now, the rules selected for the control element are certainly inadequate to control the full-scale pH system; the one that includes the changing process dynamics. However, the performance of a FLC can be dramatically altered by changing the membership functions. This is equivalent to changing the definition of the terms used to describe the variables being considered by the controller. As will be seen shortly, GAs are powerful tools capable of rapidly locating efficient fuzzy membership functions that allow the controller to accommodate changes in the dynamics of the pH system.

Analysis Element

The analysis element recognizes changes in parameters associated with the problem environment not taken into account by the rules used in the control element. In the pH system, these parameters include: (1) the concentrations of the acid and base of the input control streams, (2) the flow rates of the acid, the base, and the buffer that are randomly altered, and (3) the system setpoint. Changes to any of these parameters can dramatically alter the way in which the system pH responds to additions of acid or base, thus forming a new problem environment requiring an altered control strategy. Recall that the FLC used for the control element presented includes none of these parameters in its 14 rules. Therefore, some mechanism for altering the prescribed actions must be included in the control system. But before the control element can be altered, the control system must recognize that the problem environment has changed, and compute the nature and magnitude of the changes.

The analysis element recognizes changes in the system parameters by comparing the response of the physical system to the response of a model of the pH system. In general, recognizing changes in the parameters associated with the problem environment requires the control system to store information concerning the past performance of the problem environment. This information is most effectively acquired through either a data base or a computer model. Storing such an extensive data base can be cumbersome and requires extensive computer memory. Fortunately, the dynamics of the pH system are well understood for buffered reactions, and can be modeled using a single cubic equation that can be solved for $[\text{H}_3\text{O}^+]$ ion concentrations, to directly yield the pH of the solution. In the approach adopted here, a computer model predicts the response of the laboratory pH system. This predicted response is compared to the response of the physical system. When the two responses differ by a threshold amount over a finite period of time, the physical pH system is considered to have been altered.

When the above approach is adopted, the problem of computing the new system parameters becomes a curve fitting problem [8]. The parameters associated with the computer model produce a particular response to changes in the action variables. The parameters must be selected so that the response of the model matches the response of the actual problem environment.

An analysis element has been forged in which a GA is used to compute the values of the parameters associated with the pH system. When employing a GA in a search problem, there are basically two decisions that must be made: (1) how to code the parameters as bit strings and (2) how to evaluate the merit of each string (the fitness function must be defined). The GA used in the analysis element employs concatenated, mapped, unsigned binary coding [6]. The bit-strings produced by this coding strategy were of length 200: the first 40 bits of the strings were used to represent the concentration of the acid on the control input stream, the second 40 bits were used to represent the concentration of the base on the control input stream, the third 40 bits were used to represent the flow rate of the acid of the external streams, and the final 80 bits were used to represent the flow rates of the buffer and the base of the external streams, respectively. The 40 bits associated with each individual parameter were read as a binary number, converted to decimal numbers ($000 = 0$, $001 = 1$, $010 = 2$, $011 = 3$, etc.), and mapped between minimum and maximum values according to the following:

$$C = C_{\min} + \frac{b}{(2^m - 1)} (C_{\max} - C_{\min}) \quad (1)$$

where C is the value of the parameter in question, b is the binary value, m is the number of bits used to represent the particular parameter (40), and C_{\min} and C_{\max} are minimum and maximum values associated with each parameter that is being coded.

A fitness function has been employed that represents the quality of each bit-string; it provides a quantitative evaluation of how accurately the response of a model using the new model parameters matches the response of the actual physical system. The fitness function used in this application is:

$$f = \sum_{i=0s}^{i=100s} (pH_{model} - pH_{actual})^2. \quad (2)$$

With this definition of the fitness function, the problem becomes a minimization problem: the GA must minimize f , which as it has been defined, represents the difference between the response predicted by the model and the response of the laboratory system.

Figure 4 compares the response of the physical pH system to the response of the simulated pH system that uses the parameters determined by a GA. This figure shows that the responses of the computer model and the physical system are virtually identical, thereby demonstrating the effectiveness of a GA in this application. The GA was able to locate the correct parameters after only 500 function evaluations, where a function evaluation consisted of simulating the pH system for 100 seconds. Locating the correct parameters took approximately 20 seconds on a 386 personal computer. Industrial systems may mandate that a control action be taken in less than 20 seconds. In such cases, the time the GA is allotted to update the model parameters can be restricted. Once new parameters (and thus the new response characteristics of the problem environment) have been determined, the adaptive element must alter the control element.

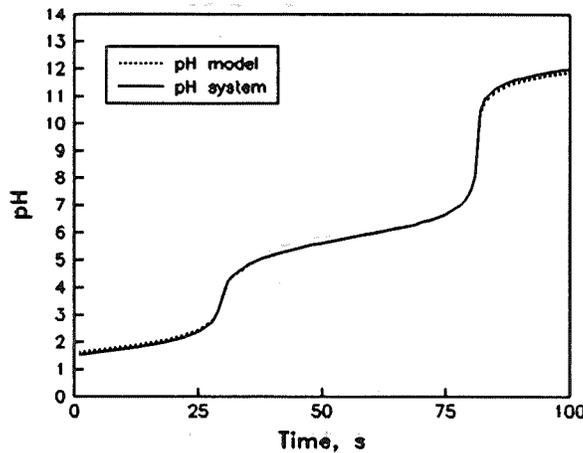


Fig. 4. Performance of an analysis element.

Learning Element

The learning element alters the control element in response to changes in the problem environment. It does so by altering the membership functions employed by the FLC of the control element. Since none of the randomly altered parameters appear in the FLC rule set, the only way to account for these conditions (outside of completely revamping the system) is to alter the membership functions employed by the FLC. These alterations consist of changing both the position and location of the trapezoids used to define the fuzzy terms.

Altering the membership functions (the definition of the fuzzy terms in the rule set) is consistent with the way humans control complex systems. Quite often, the rules-of-thumb humans use to manipulate a problem environment remain the same despite even dramatic changes to that environment; only the conditions under which the rules are applied are altered. This is basically the approach that is being taken when the fuzzy membership functions are altered.

The U.S. Bureau of Mines uses a GA to alter the membership functions associated with FLCs, and this technique has been well documented [7]. A learning element that utilizes a GA to locate high-efficiency membership functions for the dynamic pH laboratory system has been designed and implemented.

The performance of a control system that uses a GA to alter the membership functions of its control element is demonstrated for two different situations. First, Fig. 5 compares the performance of the adaptive control system (one that changes its membership functions in response to changes in the system parameters) to a non-adaptive control system (one that ignores the changes in the system parameters). In this figure, the pH system has been perturbed by the addition of an acid (at 75 seconds), a base (at 125 seconds), and a buffer (at 175 seconds). In this case, the process dynamics are dramatically altered due to the addition of the buffer, and the adaptive controller is better.

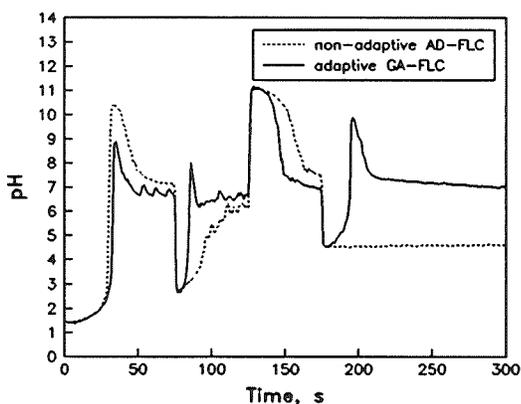


Fig. 5. External reagent additions.

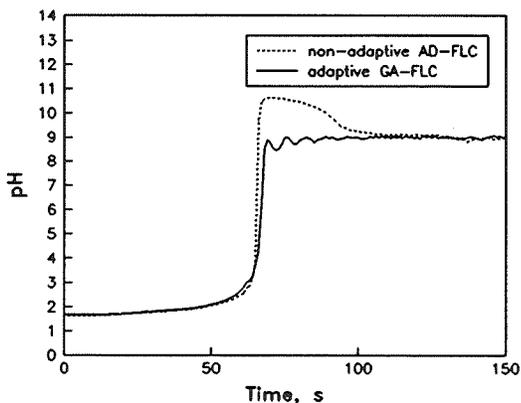


Fig. 6. Alteration of titrant concentrations.

Second, the concentrations of the acid and base the FLC uses to control pH are changed (those from the control input streams), which causes the system to respond differently. For example, if the 0.1 M HCl is the control input,

the pH falls a certain amount when this acid is added. However, all other factors being the same, the pH will not fall as much when the same volume of the 0.05 M HCl is added. The results of this situation are summarized in Fig. 6. In this simulation, the concentration of the titrants is changed at 50 seconds. As above, the adaptive control system is more efficient.

SUMMARY

Scientists at the U.S. Bureau of Mines have developed an AI-based strategy for adaptive process control. This strategy uses GAs to fashion three components necessary for a robust, comprehensive adaptive process control system: (1) a control element to manipulate the problem environment, (2) an analysis element to recognize changes in the problem environment, and (3) a learning element to adjust to changes in the problem environment. The application of this strategy to a laboratory pH system has been described.

REFERENCES

- [1] Kelly, E. G. and Spottiswood, D. J. (1982). *Introduction to Mineral Processing*. John Wiley & Sons, New York, NY.
- [2] Fogler, H. S. (1986). *Elements of Chemical Reaction Engineering*. Prentice-Hall, Englewood Cliffs, NJ.
- [3] Gottinger, W. W. (1991). *Economic Models and Applications of Solid Waste Management*. Gordon and Breach Science Publishers, New York, NY.
- [4] Zadeh, L. A. (1973). Outline of a new approach to the analysis of complex systems and decision processes. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-3, 28-44.
- [5] Goldberg, D. E. (1989). *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley, Reading, MA.
- [6] Karr, C. L. and Gentry, E. J. (1992, in press). *An Adaptive System for Process Control*. U.S. Bureau of Mines Report of Investigations.
- [7] Karr, C. L. (1991). Genetic algorithms for fuzzy logic controllers. *AI Expert*, 6, 26-33.
- [8] Karr, C. L., Stanley, D. A., and Scheiner, B. J. (1991). *A Genetic Algorithm Applied to Least Squares Curve Fitting*. U.S. Bureau of Mines Report of Investigations No. 9339.

552-61
150522

N93-25603 10

A GENETIC ALGORITHM TOOL (SPLICER) FOR COMPLEX SCHEDULING PROBLEMS AND THE SPACE STATION FREEDOM RESUPPLY PROBLEM

Lui Wang
Software Technology Branch
NASA Johnson Space Center
Houston, Texas

Manuel Valenzuela-Rendon
ITESM, Campus Monterrey
Center for Artificial Intelligence
Monterrey, N.L., Mexico

Abstract

The Space Station Freedom will require the supply of items in a regular fashion. A schedule for the delivery of these items is not easy to design due to the large span of time involved and the possibility of cancellations and changes in shuttle flights. This paper presents the basic concepts of a genetic algorithm model, and also presents the results of an effort to apply genetic algorithms to the design of propellant resupply schedules. As part of this effort, a simple simulator and an encoding by which a genetic algorithm can find near optimal schedules have been developed. Additionally, this paper proposes ways in which robust schedules, i.e. schedules that can tolerate small changes, can be found using genetic algorithms.

1. Introduction

A schedule for the delivery of resupplies to the Space Station Freedom is not easy to design due to the large span of time involved and the possibility of cancellations and changes in shuttle flights. Additionally, there is difficulty in defining ways to determine the quality of schedules, many factors should be optimized. Genetic algorithms seem adequate for the task of finding appropriate schedules due to their proven ability to deal with complex objective functions and their robustness as search methods.

An initial effort to study the application of genetic algorithms to finding schedules for the resupply of propellant to the space station has been undertaken at the Software Technology Branch of NASA Johnson Space Center. As a result of this effort, a simulator was developed. A representation of schedules was designed. Also, an evaluation mechanism was proposed.

This paper first briefly describes the basic concepts of a genetic algorithm model, followed by a short description of the design of a genetic algorithm tool (Splicer) which was developed by the Software Technology Branch. Finally, the paper describes the results of the effort to apply genetic algorithms to the design of schedules for the resupply of propellant to the Space Station Freedom.

2. Genetic Algorithms and Splicer

2.1 Genetic Algorithms

Genetic algorithms (GA) are highly parallel, mathematical, adaptive search procedures based loosely on the processes of natural genetics and Darwinian survival of the fittest. These

algorithms apply genetically-inspired operators to populations of potential solutions in an iterative fashion, creating new populations while searching for an optimal (or near-optimal) solution to the problem at hand. There are several key features of this search and optimization technique. One, the problem space is searched in parallel based on the "building blocks" concept. Two, genetic algorithms are very effective when searching function spaces that are not smooth or continuous—functions which are very difficult (or impossible) to search using calculus based methods. Three, genetic algorithms are blind: that is, they know nothing of the problem being solved other than payoff or penalty (i.e., objective function) information.

The basic iterative model of the genetic algorithms is: the algorithm starts up with an random population, and subsequent populations are created from the previous population by means of *evaluation, selection, and reproduction*. This process repeats itself until the population converges on an optimal solution or some other stopping condition is reached.

The initial population consists a set of individuals (i.e., potential solutions) generated randomly or heuristically. In the classical genetic algorithm, each member is represented by a fixed-length binary string of bits (a *chromosome*) that encodes parameters of the problem. This encoded string can be decoded to give the integer values for these parameters.

Once the initial population has been created, the evaluation phase begins. The genetic algorithms require that members of the population can be differentiated according to *goodness* or *fitness*. The members that are more fit are given a higher probability of participating during the selection and reproduction phases. Fitness is measured by decoding a chromosome and using the decoded parameters as input to the objective function. The value returned by the objective function (or some transformation of it) is used as the fitness value.

During the selection phase, the population members are given a target sampling rate which is based on fitness and determines how many times a member will mate during this generation—that is, how many offspring from this individual will be created in the next population. The target sampling rate (usually not a whole number) must be transformed into an integer number of matings for each individual. There are many ways of determining the target sampling rate and the actual number of matings. Suffice it to say that individuals that are more fit are given a reproductive advantage over less fit members.

During the reproduction phase, two members of the mating pool (i.e., members of the population with non-zero mating counts) are randomly chosen and genetic operators are applied to their genetic material to produce two new members for the next population. This process is repeated until the next population is filled. The recombination phase usually involves two operators: crossover and mutation. During crossover, the two parents exchange substring information (genetic material) at a random position in the chromosomes to produce two new strings. Crossover occurs according to a crossover probability, usually between 0.5 and 1.0. The crossover operation searches for better *building blocks* within the genetic material which combine to create optimal or near-optimal problem parameters and, therefore, problem solutions, when the string is decoded. Mutation is a secondary operation in the genetic algorithm process. It is used to maintain diversity in the population—that is, to keep the population from prematurely converging on one solution and to create genetic material that may not be present in the current population. The mechanics of the mutation operation are simple: for each position in a string created during crossover, change the value at that position according to a mutation probability. The mutation probability is usually very low—less than 0.05.

2.2 Splicer

The Splicer tool is a project within the Software Technology Branch. The purpose of the project is to develop a tool that will enable the widespread use of genetic algorithm technology.

The design chosen for the Splicer consists of four components: a genetic algorithm kernel and three types of interchangeable libraries or modules: representation libraries, fitness modules, and user interface libraries.

A *genetic algorithm kernel* was developed that is independent of representation (i.e., problem encoding), fitness function, or user interface type. The GA kernel comprises all functions necessary for the manipulation of populations. These functions include the creation of populations and population members, the iterative population model, fitness scaling, parent selection and sampling, and the generation of population statistics. In addition, miscellaneous functions are included in the kernel (e.g., random number generators). Different types of problem-encoding schemes and functions are defined and stored in *interchangeable representation libraries*. This allows the GA kernel to be used for any representation scheme. At present, the Splicer tool provides representation libraries for binary strings and for permutations. These libraries contain functions for the definition, creation, and decoding of genetic strings, as well as multiple crossover and mutation operators. Furthermore, the Splicer tool defines the appropriate interfaces to allow users to create new representation libraries (e.g., for use with vectors or grammars).

Fitness functions are defined and stored in interchangeable *fitness modules*. Fitness modules are the only piece of the Splicer system a user will normally be required to create or alter to solve a particular problem. Within a fitness module, a user can create a fitness function, set the initial values for various Splicer control parameters (e.g., population size), create a function which graphically draws the best solutions as they are found, and provide descriptive information about the problem being solved. The tool comes with several example fitness modules.

The Splicer tool provides three *user interface libraries*: a Macintosh user interface, an X Window System user interface, and a simple, menu-driven, character-based user interface. The first two user interfaces are event-driven and provide graphic output using windows.

The C programming language was chosen for portability and speed. Splicer has been tested on multiple platforms which include Sun 3/80™, SPARC™, IBM RS6000™, and Apple Macintosh™. With the new character and TTY interfaces, Splicer can now be embedded in the user application.

3. The Space Station Freedom scheduling problem

The Space Station Freedom will require the supply of various items in a regular fashion, including such things as air, food, experiment payload modules, and propellant. A schedule for the delivery of these items spanning several years will be needed. Because of the large number of activities to be scheduled, and the long period of time involved, it is not possible to plan them all at the same time in a single schedule. Instead, independent schedules for separate items will be designed. Then, these will be merged into a single schedule. To facilitate this integration, each individual schedule must be as flexible and robust as possible, i.e. changes must be easily made and must not seriously affect the overall performance. Additionally, the overall schedule must also be robust so that shuttle flight delays and cancellations can be tolerated without major changes. Genetic Algorithms (Goldberg, 1989, Holland, 1975), due to their inherent flexibility, seem to be an appropriate tool for solving this problem because the complexity of the many restrictions that can be involved in schedules for individual items.

As a first step in studying the applicability of genetic algorithms, the problem of scheduling the resupply of propellant was selected. The following description of propellant resupply to the Space Station Freedom was a result of several talks with the Level II Space Station Freedom Resource Utilization Analysis engineers.

3.1 Description of the propellant resupply problem

Reboosting at the space station occurs after every departure of the space shuttle. The reboost operation takes the space station to its highest orbit. Between reboosts, the space station slowly loses altitude, so as to meet the shuttle at its lowest orbit.

The thrust required for a reboost operation is supplied by three out of six reboosting modules. These modules contain a propellant which is consumed to produce a force on the space station. The space station has eight parking spaces where reboosting modules can be placed. These eight spaces are grouped in pairs and each pair is located on a corner of an imaginary rectangle perpendicular to the Earth. There are a total of eight modules; at any given time, six modules are expected to be on the space station and two on ground.

For any reboost operation to be carried out, three modules in different corners must be fired; in this way, a force and a torque, in any desired direction perpendicular to the Earth, can be applied to the space station. For a given reboost, and assuming that the modules contain sufficient propellant, there will always be two sets of three corners of the rectangle to choose from, i.e. two triangles (call them A and B) that can produce the same force and torque. Both triangles will be able to produce the same reboosting effect, and require the same total propellant, but will consume different levels of propellant from the individual corners involved. The level of propellant required for a given reboost operation is not constant; it depends on the solar activity at the moment, and thus, can only be estimated ahead of time.

The space station will normally have five reboost cycles per year, one for each flight of the space shuttle. At its lowest, orbit it will meet the space shuttle and receive resupplies, then it will be reboosted to its highest orbit. The resupply of propellant will normally consist of delivering two full modules to the space station and removing two whose propellant level is either very low or zero. The removed modules will be returned to Earth and refurbished; this operation will take a lapse of time equal to the time between two shuttle flights, and thus, returned modules must stay on Earth for at least one shuttle flight. During these resupply operations, modules that are not empty can be moved from a parking space to another. Considering the levels of propellant required for typical reboost operations, Level II Space Station Freedom Resource Utilization Analysis Office estimate that there will be approximately a delivery of two modules every year. Besides reboosting, propellant is needed for *attitude control*. The requirement for attitude control is estimated as a fixed quantity equally distributed among the four corners of the rectangle where modules are parked.

The space station must have enough propellant to continue operation in spite of several contingencies. An operation called *collision avoidance maneuver* requires fixed (and unequal) quantities of propellant from the four corner of the rectangle. Additionally, normal reboosting should be possible if a scheduled shuttle flight is canceled, i.e. the space station should be able to perform a *skip cycle* reboost even if the canceled shuttle flight is one in which propellant was to be supplied. The requirements for a skip cycle are those of a reboost and attitude control. It must be underlined, that these contingency propellant requirements are not actually consumed, but must be available.

3.2 Propellant resupply schedule form

A schedule for all the operations regarding propellant is needed. The schedule should answer four basic questions:

- What triangle is to be fired on each reboost operation?
- On which shuttle flights should propellant modules be supplied to the space station?
- For each propellant resupply operation, which modules should be removed, and where should the new modules be placed?

- For each propellant resupply operation, which modules, if any, should be moved? and to which corners of the rectangle?

3.3 Optimization goals

The complete criteria for defining what is a good schedule has not been fully determined yet. The following elements are known to be desired in a good schedule.

- **Robustness:**
A good schedule should be capable of tolerating small changes due to modifications or cancellations in shuttle flights, variations in solar activity, requirements of schedules for the supply of other items, etc.
- **Wasted propellant:**
Propellant contained in modules returned to Earth is a waste and should be minimized. A efficient schedule should ask for the return of modules only when these are empty or almost empty.
- **Propellant resupply operations:**
A good schedule will require the minimal number of propellant resupply flights.
- **Module movements on the space station:**
As explained before, non empty modules can be moved from one to corner to another during a resupply operation. A good schedule should make the least number of these movements.

3.4 Simulator outline

A simple simulator, written in C and running on the Macintosh and on UNIX workstations, was developed for this problem. The simulator allows a user to type in instructions regarding the use of propellant and its resupply. For each shuttle flight, the user is asked if a resupply should occur; if so, it is then asked if modules should be moved from one corner to another, which modules should be returned to Earth and where should the new modules be placed. The user is also asked which triangle should be fired on each reboost operation. The simulator displays the placement of each modules and its propellant contents, and takes into account the propellant requirements for contingencies. The propellant content of the modules is displayed to the user after each operation. In this way, a whole schedule can be tested step by step.

In summary, the user interacts with the simulator by making decisions regarding the following operations:

- **Reboost:**
The user decides which triangle of propellant modules is to be fired on each reboost cycle. The simulator displays the contents of the modules. In case of a schedule failure, i.e. the user has drained modules so that a reboost is impossible, the simulator traces back in time to the last viable situation so that the user can continue simulation from this point.
- **Resupply of propellant:**
The simulator allows the user to decide when to perform a resupply, which modules to return, and where to place the new modules.
- **Movements of modules:**
If a resupply is to occur, the simulator asks the user which modules are to be moved, and which are the target corners. The user can enter as many movements as desired.

3.5 Simulator reinterpretation of commands

An important characteristic of the simulator, which is necessary if it is to be coupled to a genetic algorithm, is its capacity to reinterpret commands when they cannot be carried out.

When the user requests an operation that is not valid, the simulator is capable of either ignoring the request or finding the next possible operation of the same type which is possible. The following is a list of the reinterpretations the simulator performs:

- **Try other triangle**
When the user attempts to perform a reboost using a triangle whose modules do not contain enough propellant, the simulator tests the other triangle, and if valid, it fires it.
- **Two consecutive resupply operations**
Because of the requirement that modules that return to Earth must stay for at least one shuttle flight, resupply operations cannot be consecutive. The simulator ignores any attempts to perform consecutive resupplies.
- **Removing empty modules**
During a resupply operation, empty modules must be removed. The simulator ignores attempts to remove modules that are not empty if other that are empty are not removed first.
- **Removing least full module of corner**
During a resupply operation, the least full module of a corner must be removed. When specifying modules to be removed, the user only indicates the desired corner; the simulator removes the least full module from that corner.
- **Placing modules on corners that are full**
Each corner of the rectangle contains two parking spaces for modules. If the users tries to place a module on a corner that already has two modules, the simulator will place it on the next corner that has an empty parking space.

These reinterpretations of user commands are important when the simulator is coupled to a genetic algorithm, because in this way many invalid individuals are expressed as valid schedules, and thus, are efficiently employed in the genetic search. In effect, through these transformations, the search space is pruned of regions that contain only invalid solutions.

4. Coupling the simulator to Splicer

There are two main issues to resolve when coupling a simulator, like the one developed for this work, and a genetic algorithm tool such as Splicer. First, a way to encode possible solutions into binary strings must be chosen. Second, a method to evaluate this solutions must be defined. In the following section these matters are explained.

4.1 Encoding of schedules into binary strings

The possible schedules must be encoded into binary strings, or chromosomes, so that the genetic algorithm can operate over a population of them. Thus, an encoding is a transformation from the space of possible schedules, i.e. lists of commands that can be given to the simulator, to the space of binary strings. In this encoding design the chromosome is composed of four segments and each segment encodes one of the possible commands to the simulator.

- **Triangle selection**
This segment consists of as many bits as reboosts being scheduled. Each position indicates with a 0 that triangle A, or with a 1 that triangle B, is to be fired at the corresponding reboost.
- **Is resupply to be performed?**
This segment consists of as many bits as reboosts being scheduled. Each position indicates if a resupply is to be performed.
- **Selection of corners for resupply**

This segment consists of as many subsegments as necessary to represent the maximum number of resupplies thought necessary for the space station to operate for the lapse of time being scheduled. Each subsegment contains eight bits in groups of two. Every two bits point to one of the four corners.

- Selection of modules to be moved

This segment consists of as many subsegments as necessary to represent the maximum number of module movements thought necessary. Each subsegment is composed of four bits. The first two bits indicate which module is to be moved and the last two indicate the corner to where it will be moved. Notice that during a resupply, two modules must be chosen to be removed so that there are only four left to be moved to different corners.

Each chromosome, or individual in genetic algorithm, encodes the operations a user would request from the simulator. Even though it can be handled by the simulator that was developed, the selection of modules to be moved was not incorporated into the genetic search. This step remains for future work.

4.2 Evaluation

The problem of finding appropriate schedules for the resupply of propellant can be considered at two levels. First, and most important, it is necessary to find *valid* schedules, i.e. schedules that can perform all the reboosts asked for. Then, minimization of wasted propellant and maximization of robustness can be considered. In a genetic algorithm, where initial individuals are generated randomly, it is very unlikely that valid solutions will be contained in the initial population. Thus, the first goal of the evaluation procedure should be to drive the population to valid solutions. In light of the above, the objective function chosen had the following form:

$$\text{fitness} = \text{number of successful flights} \cdot 10 \\ - \text{total wasted propellant} / 1000.$$

The maximization of robustness has yet to be considered. The following section describes ways under consideration to include robustness in the optimization.

5. Maximization of robustness and multi-objective optimization

The objective of maximizing robustness in a schedule requires a different approach than that used to attack all other optimization goals. Robustness is difficult to measure. In the context of resupply to the space station, a robust schedule is one that can tolerate small changes due to modification or cancellations in shuttle flights, variations expected in solar activity, requirements of schedules for the supply of other items, etc. What is meant by "tolerate small changes" is yet to be defined. Additionally, robustness is a characteristic that will surely conflict with other optimization goals; for example, a schedule that requires for propellant resupply more often than necessary will probably be more tolerant to shuttle flight cancellations, but will waste propellant.

The problem of maximizing robustness can be addressed in two steps. First, develop a stochastic method to measure robustness of a schedule and incorporate it into the objective function. Second, the problem will be treated as one of multi-objective optimization where robustness will be maximized independent by of all other objectives. In the following subsections the proposing method for applying the genetic algorithm is presented.

5.1 Measuring robustness

Genetic algorithms are known to be tolerant to noise in the evaluation of the objective function (Fitzpatrick, Grefenstette, and Van Gucht, 1984). The reason for this is that the genetic

algorithm implicitly processes *schemata* which are patterns of bits that are represented by many members of the population. In this way, the evaluation of a single individual can be corrupted with noise, and the genetic algorithm will continue to find near optimal solutions, as long as the noise has zero mean. Therefore, for each individual schedule in the population, a number of random small changes will be generated. The robustness will be obtained as the average performance of the modified schedules; in the expected value, this will be a correct measure of robustness. This value of robustness can be incorporated into the objective function as part of a linear combination of all the optimization goals. In this way, for a given relation of relative importance among the optimization goals, the genetic algorithm will find a near optimal solution.

5.2 Multi-objective optimization with genetic algorithms

Many real optimization problems, like the scheduling of supplies to the space station, require that several criteria be optimized simultaneously. In these problems, some of the objectives are conflicting, and it is difficult to decide the relative importance of each one. In practice, all the criteria are usually combined into a single objective function by taking a linear combination of them, in order to avoid the problem of multi-objective optimization. This is not always the best approach, but one often taken because of limitations on the optimization methods.

The area of multi-objective optimization has not been fully attacked by genetic algorithm researchers. Schaffer (1985) has presented a method for applying genetic algorithms to this type of problem, and no continuation of this effort has been reported. Goldberg (1989, p. 201) suggested applying a combination of rank selection (Baker, 1985) and niche and speciation methods (Deb and Goldberg, 1989). An extension to the previous efforts, (Valenzuela Rendón & Cantú Aguillen, 1992) that apply niche and speciation methods to the solution of multimodal problems to include multi-objective optimization as proposed by Goldberg, is proposed to attack this problem. In this way, the genetic algorithm will find not one but a family of near optimal schedules for varying degrees of relative importance of robustness versus all other optimization goals.

6. Final Comments

The Space Station Freedom will require the resupply of many items over a long period of time. Schedules for individual items will first be designed, and thus, these schedules will be integrated into a single schedule. For this integration to take place, individual schedules must be able to tolerate small changes. Genetic algorithms seem like an appropriate tool to apply to this problem in view of their proven ability to solve complex objective functions and their robustness as search methods.

As a first step, we have developed a simple simulator and an encoding by which a genetic algorithm can be applied to finding schedules for the resupply and use of propellant without considering robustness. As following steps, we will take advantage of the ability of genetic algorithms to tolerate noise in the objective function, and measure robustness in a stochastic manner. Also, we propose to apply multi-objective optimization techniques in genetic algorithms to find families of near optimal schedules for varying degrees of relative importance of robustness.

Acknowledgments

This work was performed during a six month stay of Manuel Valenzuela-Rendón at NASA's Software Technology Branch as part of RICIS (Research Institute for Computer and Information Systems) research activity number SR.04 (NASA Cooperative Agreement NCC-9-16) and with partial support of the ITESM.

References

- Baker, J. E. (1985). Adaptive selection methods for genetic algorithms. *Proceedings of an International Conference on Genetic Algorithms and Their Applications*, pp. 101–111.
- Etter, D. M., Hicks, M. J., & Cho, K. H. (1982). Adaptive genetic algorithm for determining optimum filter coefficients in recursive adaptive filter. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, Vol. 2, pp. 635–638.
- Davis, L. (1985). Applying adaptive algorithms to epistatic domains. *Proceedings of an International Conference on Genetic Algorithms and Their Applications*, pp. 162–164.
- Deb, K., & Goldberg, D. E. (1989). An investigation of niche and species formation in genetic function optimization. *Proceedings of the Third International Conference on Genetic Algorithms*, pp 42–50.
- Fitpatrick, J. M., Grefenstette, J. J., & Van Gucht, D. (1984). Image registration by genetic search. *Proceedings of IEEE Southeast Conference*, 460–464.
- Goldberg, D. E. (1983). Computer-aided gas pipeline operation using genetic algorithms and rule learning (Doctoral dissertation, University of Michigan). *Dissertation Abstracts International*, 44(10), 3174B. (University Microfilms No. 8402282).
- Goldberg, D. E. (1989). *Genetic algorithms in search, optimization, and machine learning*. Reading, MA: Addison-Wesley.
- Goldberg, D. E. & Lingle, R. (1985). Alleles, loci, and the traveling salesman. *Proceedings of an International Conference on Genetic Algorithms and Their Applications*, pp. 154–159.
- González Sáenz, A. (1991). *Algoritmos genéticos: Una aplicación al diseño de filtros digitales* [Genetic algorithms: An application to digital filter design]. Unpublished master's thesis, ITESM, Monterrey.
- Guerra-Salcedo, C. M., & Valenzuela-Rendón, M. (1991). Resolviendo consultas de apareamiento parcial utilizando algoritmos genéticos [Solving partial match queries by means of genetic algorithms]. *Memorias de la VIII Reunión Nacional de Inteligencia Artificial*, 13–28. Sociedad Mexicana de Inteligencia Artificial.
- Holland, J. H. (1975). *Adaptation in natural and artificial systems*. Ann Arbor, MI: University of Michigan Press.
- Shaffer, J. D. (1985). Multi-objective optimization with vector evaluated genetic algorithms. *Proceedings of an International Conference on Genetic Algorithms and Their Applications*, pp. 93–100.
- Valenzuela-Rendón, M., Guerra-Salcedo, C. M., & Icaza, J. I. (1991). A genetic algorithm approach to partial match retrieval based on hash functions. *Proceedings of the IV International Symposium on Artificial Intelligence*, 156–162.

Valenzuela Rendón, M. & Cantú Aguillen, C. (1992). *Solución a sistemas de ecuaciones no lineales simultáneas usando algoritmos genéticos*. [Solution to nonlinear simultaneous equations by genetic algorithms]. Manuscript submitted for publication.

Wang L. (1991). Genetic Algorithm Overview, *proceedings of The 1991 Science, Engineering & Technology Seminars*

02/11/17

**ENERGY AND ENVIRONMENT PART 3:
ENVIRONMENTAL TECHNOLOGIES**



**"Preliminary Studies Leading Toward the Development
of a LIDAR Bathymetry Mapping Instrument"**

Dr. John M. Hill, Mr. Brendan D. Krenek, and Dr. Terry D. Kunz;
Space Technology And Research (STAR) Center, Houston Advanced Research Center
(HARC), The Woodlands, Texas

Mr. William Krabill and Mr. Fran Stetina;
NASA/Goddard Space Flight Center, Greenbelt, MD

ABSTRACT:

The National Aeronautics and Space Administration (NASA) at Goddard Space Flight Center (GSFC) has developed a laser ranging device (LIDAR) which provides accurate and timely data of earth features. NASA/GSFC recently modified the sensor to include a scanning capability to produce LIDAR swaths. They have also integrated a Global Positioning System (GPS) and an Inertial Navigation System (INS) to accurately determine the absolute aircraft location and aircraft attitude (pitch, yaw, and roll), respectively. The sensor has been flown in research mode by NASA for many years. The LIDAR has been used in different configurations or modes to acquire such data as altimetry (topography), bathymetry (water depth), laser-induced fluorosensing (tracer dye movements, oil spills and oil thickness, chlorophyll and plant stress identification), forestry, and wetland discrimination studies.

NASA and HARC are developing a commercial version of the instrument for topographic mapping applications. The next phase of the commercialization project will be to investigate other applications such as wetlands mapping and coastal bathymetry. In this paper we report on preliminary laboratory measurements to determine the feasibility of making accurate depth measurements in relatively shallow water (approximately 2 to 6 feet deep) using a LIDAR system. The LIDAR bathymetry measurements are relatively simple in theory. The water depth is determined by measuring the time interval between the water surface reflection and the bottom surface reflection signals. Depth is then calculated by dividing by the index of refraction of water. However, the measurements are somewhat complicated due to the convolution of the water surface return signal with the bottom surface return signal. Therefore in addition to the laboratory experiments, computer simulations of the data were made to show these convolution effects in the return pulse waveform due to: a) water depth, and b) changes in bottom surface reflectivity.

Biographical Information

- J. Hill:** Jack Hill is the Head of the Remote Sensing/Geographic Information Systems (RS/GIS) Laboratory at the Houston Advanced Research Center (HARC). He is the Principal Investigator at HARC on the NASA LIDAR Commercialization Program.
- B. Krenek:** Brendan Krenek is a Laser Systems Specialist at HARC where he assists in laboratory feasibility demonstration experiments and system prototyping.
- T. Kunz:** Terry Kunz is Head of Laser Applications in the Space Technology And Research (STAR) Center at HARC. He is the Principal Scientist on the NASA LIDAR Commercialization Program.
- W. Krabill:** Bill Krabill is the Project Scientist for the NASA LIDAR Commercialization Project and is responsible for mission planning, instrument development, data analysis, and data interpretation. He is stationed at NASA Wallops Space Flight Center.
- F. Stetina:** Fran Stetina, stationed at the International Data Systems Office, Goddard Space Flight Center, is the Project Manager for NASA's LIDAR Commercialization Program.

PRELIMINARY STUDIES LEADING TOWARD THE DEVELOPMENT OF A LIDAR BATHYMETRY MAPPING INSTRUMENT

1.0 INTRODUCTION

The National Aeronautics and Space Administration (NASA) at Goddard Space Flight Center (GSFC) has developed a laser ranging device (LIDAR) which provides accurate and timely data of earth features. NASA/GSFC recently modified the sensor to include a scanning capability to produce LIDAR swaths. They have also integrated a Global Positioning System (GPS) and an Inertial Navigation System (INS) to accurately determine the absolute aircraft location and aircraft attitude (pitch, yaw, and roll), respectively. The sensor has been flown in a research mode by NASA for many years. It has been used in different configurations or modes of the LIDAR to acquire such data as altimetry (topography) [1], bathymetry [2], laser-induced fluorosensing (tracer dye movements [3], oil spills and oil thickness [4 and 5], chlorophyll and plant stress identification [6]), forestry [7], and wetland discrimination studies [8 through 11].

NASA and HARC are developing a commercial version of the instrument for topographic mapping applications. The next phase of the commercialization project will be to investigate other applications such as wetlands mapping and coastal bathymetry.

In this paper we report on preliminary laboratory measurements to determine the feasibility of making accurate depth measurements in relatively shallow water (approximately 2 to 6 feet deep) using a LIDAR system. The LIDAR bathymetry measurements are relatively simple in theory. The water depth is determined by measuring the time interval between the water surface reflection and the bottom surface reflection signals. Depth is then calculated by dividing by the index of refraction of water. However, the measurements are somewhat complicated due to the *convolution* of the water surface return signal with the bottom surface return signal. Therefore, in addition to the laboratory experiments, computer simulations of the data were made to show these convolution effects in the return pulse waveform due to: a) water depth, and b) changes in bottom surface reflectivity.

The experiment was simplified by eliminating the water column and substituting graphite plates for the top and bottom surfaces. Depth and signal return intensities were adjusted by simply moving the plates relative to each other and intercepting different amounts of the laser beam, respectively. Section 2 describes the laboratory apparatus used in these measurements. Section 3 presents the experimental results and the computer simulations. Section 4 makes general conclusions regarding the feasibility of performing LIDAR bathymetry measurements in shallow water.

2.0 LABORATORY LIDAR APPARATUS DESCRIPTION

These measurements were performed at HARC's Laser Applications Laboratory. Figure 1 shows the apparatus used for performing the LIDAR experiments. This apparatus consisted of a pulsed laser, telescope, fast photodiode/amplifier, and digital oscilloscope. The $\lambda = 532$ nm output of the Nd:YAG laser (Spectra Physics, Model GCR-11-3) emitting $E \geq 155$ mJ/pulse in a 6-7 ns pulse width was directed approximately 50 feet down the laboratory and onto two carbon plates. The top carbon plate functioned as the water surface and the bottom carbon plate served as the bottom surface. The reflected light from the top and bottom surfaces was collected by a 1" aperture telescope. A fast photodiode (Electro-Optics Technology, Model ET-2000, 200 ps rise time) and amplifier (Stanford Research Systems, Inc., Model SR440, DC-300MHz) and digital oscilloscope (LeCroy, Model 9400, 125 MHz) detected and recorded the return laser pulse waveform, respectively. Adjusting the distance between the two carbon plates simulated different water depths; moving the plates in and out of the laser beam adjusted the strength of the top and bottom surface return signals. The distances between the two plates for these experiments were 1 foot, 3 feet, and 5 feet. For each experiment, the digital scope averaged 50 waveforms before transferring the result to the plotter.

Figure 2 shows the instrumental response of the LIDAR apparatus shown in Figure 1. In this experiment, the top carbon plate was removed and the Nd:YAG laser was in *long pulse mode* (6-7 ns FWHM pulsewidth according to manufacturer's specifications). As seen in Figure 2, the LIDAR system is recording the laser FWHM as approximately 12.8 ns. This is attributed to two factors:

1. The bandwidth of the digital scope is 125 MHz; therefore, its rise time is approximately 2.8 ns (*i.e.*, $.35/125$ MHz), and
2. The bandwidth of the amplifier is 300 MHz, corresponding to a rise time of 1.2 ns.

The Spectra Physics Nd:YAG laser can also be operated in a *short pulse mode* which reduces the temporal pulse width by approximately a factor of 3 to about 2.5 ns; this results in only a 10% loss of pulse energy. Figure 3 shows the instrumental response with the laser in the short pulse mode. After consulting with the manufacturer, the 2.5 ns specification (alignment sensitive) applies only to the base of the center peak. Therefore, in all subsequent measurements, the laser was operated in long pulse mode.

3.0 RESULTS

3.1 EXPERIMENTAL DATA and SIMULATIONS

Figures 4, 5, and 6 show the experimental results (rough curves) when the plates were separated by 1 foot, 3 feet, and 5 feet, respectively. Also shown in these figures are simulated data (smooth curves) of each experiment. The simulations were generated by convoluting two Gaussian lineshapes, each defined by a 12.8 ns FWHM, and different line intensities (as noted in the figure captions). These simulated data were generated using a program called COCON1.BAS (described below) and are simply best fits to the experimental data as determined by eye.

3.2 ADDITIONAL COMPUTER SIMULATIONS

After understanding the experimental results above, additional computer simulations were performed to reconstruct idealized bathymetry waveform data (*i.e.*, no noise due to intermediate scattering layers located between the top and bottom surfaces) predicted in 2 to 15 feet water depths. These simulations are based on the model shown in Figure 7. In this simple model, a laser pulse of intensity I_0 is shot towards the water surface. At the air - water interface a portion (I_1) of the incident beam is reflected back towards the detector based on the water index of refraction [*i.e.*, $I_1 = I_0(n-1)^2/(n+1)^2$]. The rest of the laser pulse, $I_2 (= 1-I_1)$, is transmitted down through the water column to the bottom surface. In this process, the laser is attenuated according to the Beer-Lambert law [$\log(I_{\text{final}}/I_{\text{incident}}) = -\alpha\ell$, where α is the attenuation coefficient and ℓ is the path length] and is represented as I_3 . The beam is then reflected off the bottom surface (R_{BS}) resulting in I_4 , and attenuated again travelling up the water column to the surface. A portion of I_5 is again lost at the water - air interface. I_6 is reflected back in the water, and I_7 is transmitted back towards the detector.

As shown in Figure 7, the two most important quantities for the generating reflected waveform simulation data are I_1 and I_7 . Using $n = 1.33$ as the water index of refraction at 20°C, I_1 is 2% of I_0 . In general, I_7 can be calculated by:

$$I_7 = (10^{-\alpha\ell})^2 (I_2)^2 (R_{\text{BS}}) (I_0).$$

Using this model, two families of simulations were generated. The first simulation used a water attenuation coefficient (α) of 0.1/foot [12] and bottom surface reflectivity (R_{BS}) of 40%. Reflected waveform data were simulated for water depths of 1.5, 3.0, 4.5, 6.0, 7.5, and 11.3 feet corresponding to 4, 8, 12, 16, 20, and 30 ns top and bottom surface peak separations, respectively (*i.e.*, the speed of light in water was taken to be 1.33 times less than in air). This family of curves is shown in Figure 8. An analogous family of curves was generated using a bottom surface reflectivity (R_{BS}) of 10%. These data are shown in Figure 9. As noted from Reference 12, these are reasonable values for α and R_{BS} . Each of the panels in Figures 8 and 9 list the water depth, separation time between the water surface and bottom surface return pulses, and the relative line intensities for the water surface return signal (I_{Int_1}) and the bottom surface return signal (I_{Int_2}). The data were simulated by convoluting two 6 ns FWHM Gaussian lineshapes defined by intensities

Int_1 and Int_2 . This simple convolution software, named "COCON1", was developed and run on an IBM compatible computer and written in basic. The software will convolute any number of lines using a Lorentzian or Gaussian lineshape (defined by the lineshape FWHM). Spectral resolution can be controlled by a combination of the total number of points in the file and the wavelength range of convolution. The convoluted spectrum can be written in wavenumbers (energy) or angstroms (wavelength). The results were plotted on a HP7475A using software developed in-house and named "TRANSPLT", also written in basic.

In the simulations shown in Figure 8 ($\alpha = 0.1/\text{foot}$ and $R_{BS} = 40\%$), the bottom return pulse is *stronger* than the surface return pulse for the 3 feet (8 ns peak separation) and 4.5 feet (12 ns peak separation) water depths. At water depths deeper than approximately 5 feet, the laser will be attenuated in a sufficient length of water column so that the surface return pulse will be the larger of the signals. In contrast, in the simulations shown in Figure 9 ($\alpha = 0.1/\text{foot}$ and $R_{BS} = 10\%$), the surface return signal is always stronger than the bottom return signal.

4.0 CONCLUSIONS

The following conclusions are made based on the results of these experiments and computer simulations.

1. The 125 MHz digital oscilloscope (2.8 ns rise time) was not fast enough to resolve the 6 ns FWHM pulse width of the Nd:YAG laser. Therefore, the experimental data in Figures 2, 4, 5, and 6 are not in true physical agreement with the measured block separation distances.
2. The simulated data can easily be made to agree quantitatively with the experimental data.
3. Even with the slow digital oscilloscope, top and bottom surface signals corresponding to approximate water depths of 3 feet and 5 feet were clearly discernable.
4. When the water depth is of the order of 1 to 2 feet, the recorded lineshape will simply look like a single broadened line as shown in the experimental data of Figure 4, and the top panels of the simulated data in Figures 8 and 9.
5. The simulated data in Figures 8 and 9 show the effects of the water depth and bottom surface reflectivity on the return waveform. Clearly, a trade-off exists between signal strength and signal separation between the top and bottom signal reflections. In shallow water, both top and bottom return signals should be strong, but close together in time, leading to spectrum resembling a single broadened line at 2 feet depths. As the water gets deeper, the bottom surface return pulse becomes weaker due to attenuation, but is separated farther in time from the top surface return pulse.
6. A Marquardt algorithm/fitting technique [13], or similar data reduction routine(s), would be useful for deconvoluting these lineshapes in a real bathymetry application.
7. Two (2) new digital scopes will soon be available that would nicely fill the waveform recording requirements of this bathymetry application. LeCroy will soon market their 7200/7200A digital scope equipped with a 2 Gsample/sec (500 MHz analog system bandwidth, 0.7 ns rise time) for approximately \$45,000. HP will also be marketing their 54710/54720 Mainframe digital scope equipped with a 4 Gsample/sec (1.1 GHz analog system bandwidth, 0.3 ns rise time) for approximately \$40,000. Both scopes will have a SCSI-2 (small computer systems interface) port option.
8. **Future Directions** will investigate three main topics:
 - a) range biasing as a function of scan angle,
 - b) range biasing as a function of water turbidity, and
 - c) maximum depth capability of the technique.

5.0 REFERENCES

1. W. B. Krabill, J. G. Collins, L. E. Link, R. N. Swift, M. L. Butler, "Airborne Laser Topographic Mapping Results from Joint NASA/U. S. Army Corps Of Engineers Experiments", *Photogrammetric Engineering and Remote Sensing*, 50(6), 685 (1984).
2. F. E. Hoge, R. N. Swift, and E. B. Frederick, "Water Depth Measurement Using An Airborne Pulsed Neon Laser System", *Appl. Opt.*, 19, 871 (1980).
3. W. B. Krabill and R. N. Swift, "Airborne LIDAR Experiments at the Savannah River Plant", NASA Technical Memorandum 4007, 1987.
4. F. E. Hoge and R. N. Swift, "Experimental Feasibility of the Airborne Measurement of Absolute Oil Fluorescence Spectral Conversion Efficiency", *Appl. Opt.*, 22, 37 (1983).
5. F. E. Hoge and R. N. Swift, "Oil Thickness Measurements Using Airborne Laser-Induced Water Raman Backscatter", *Appl. Opt.*, 19(19), 3269 (1980).
6. F. E. Hoge, R. N. Swift, and J. K. Yungel, "Feasibility of Airborne Detection of Laser-Induced Fluorescence of Green Plants. 1: A Technique for the Remote Detection of Plant Stress and Species Differentiation", *Appl. Opt.*, 23, 134 (1983).
7. G. A. McClean and G. L. Martin, "Merchantable Timber Volume Estimation Using An Airborne LIDAR System", *Canadian Journal of Remote Sensing*, 12(1), 7 (1986).
8. V. Carter, "Applications of Remote Sensing to Wetlands" in G. J. Johannsen and J. L. Sanders (Eds.) Remote Sensing for Resources Management, Ankeny, Iowa, Soil Conservation Society of America, 284-300 (1982).
9. M. K. Butera, "Remote Sensing of Wetlands", *IEEE Transactions on Geoscience and Remote Sensing*, GE-23, 383 (1983).
10. J. R. Jensen, M. E. Hodgson, E. J. Christensen, H. E. Mackey, and L. Tinney, "Remote Sensing of Inland Wetlands: A Multispectral Approach", *Photogrammetric Engineering and Remote Sensing*, 52, 87 (1986).
11. J. R. Jensen, E. W. Ramsey, H. E. Mackey, E. J. Christensen, and R. R. Sharitz, "Inland Wetland Change Detection Using Aircraft MSS Data", *Photogrammetric Engineering and Remote Sensing*, 52, 521 (1987).
12. F. E. Hoge, C. Wayne Wright, William B. Krabill, Rodney R. Buntzen, Gary D. Gilbert, Robert N. Swift, James K. Yungel, and Richard E. Berry, *Applied Optics*, 27(19), 3969 (1988).
13. W. Schreiner, M. Kramer, S. Krischer, and Y. Langsam, *PC Tech. J.*, 3, 170 (1985).

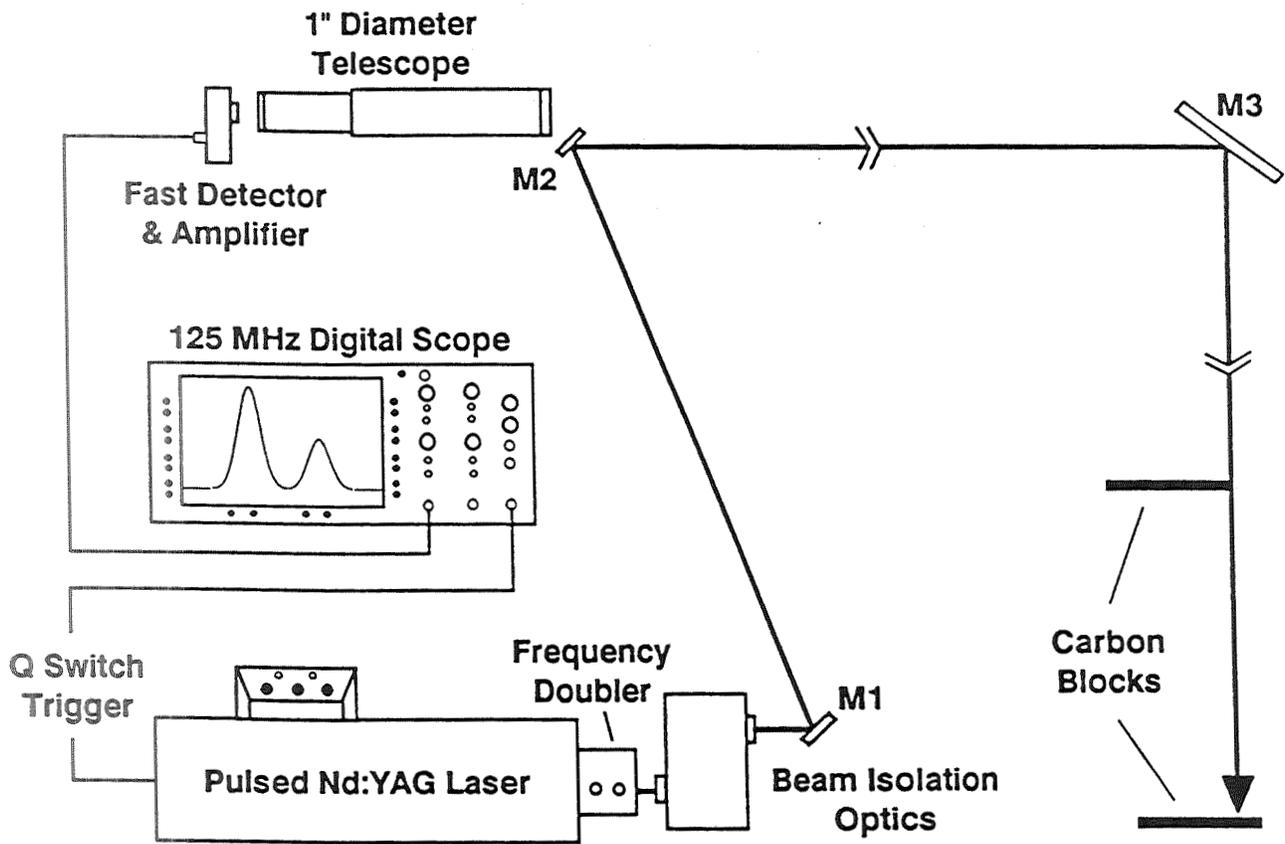


Figure 1. Laboratory LIDAR apparatus for simulating bathymetry measurements.

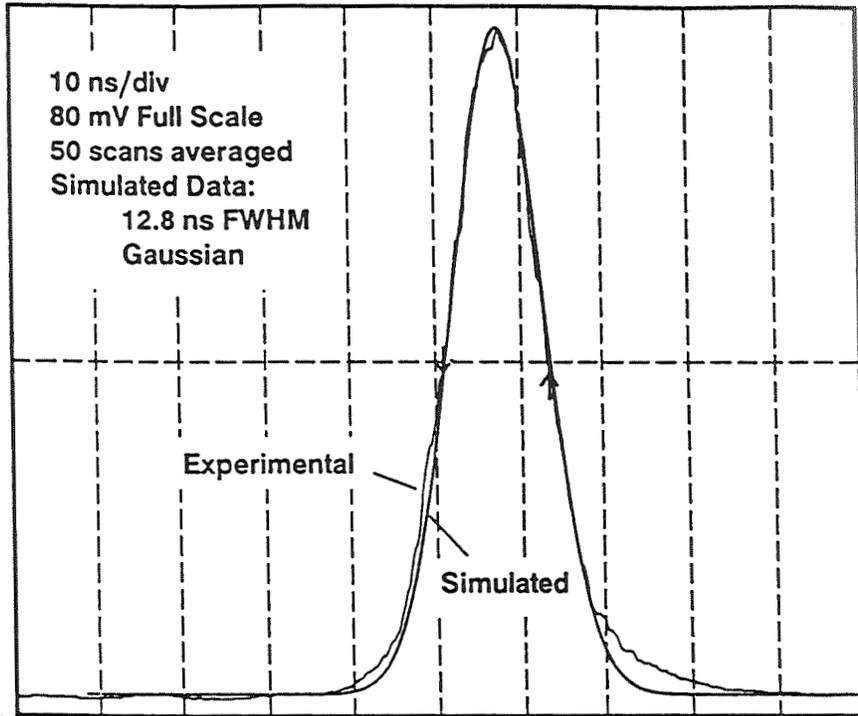


Figure 2. Instrumental response of the LIDAR apparatus shown in Figure 1 with the Nd:YAG laser in long pulse mode.

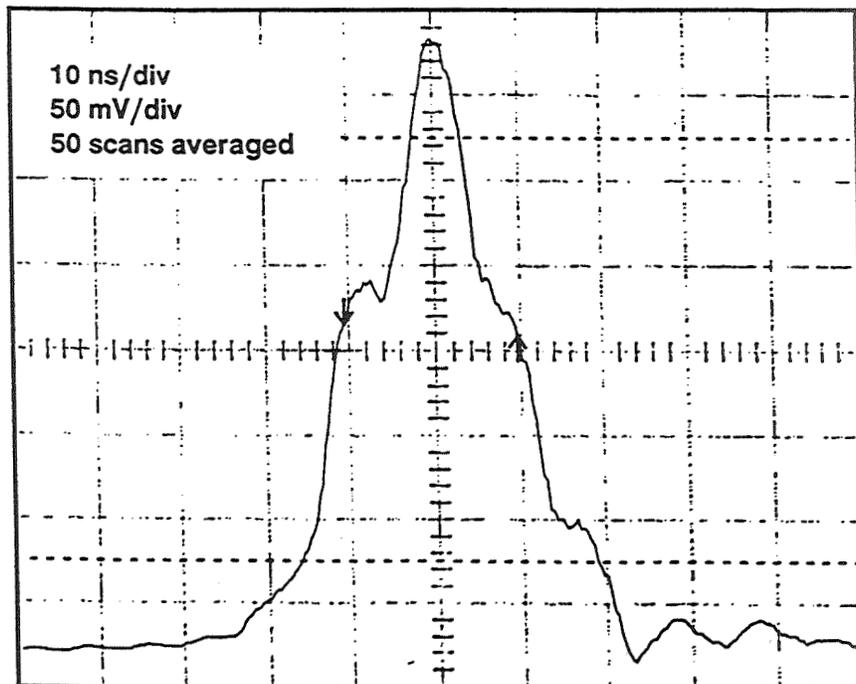


Figure 3. Instrumental response of the LIDAR apparatus shown in Figure 1 with the Nd:YAG laser in short pulse mode. Note the extra side lobes.

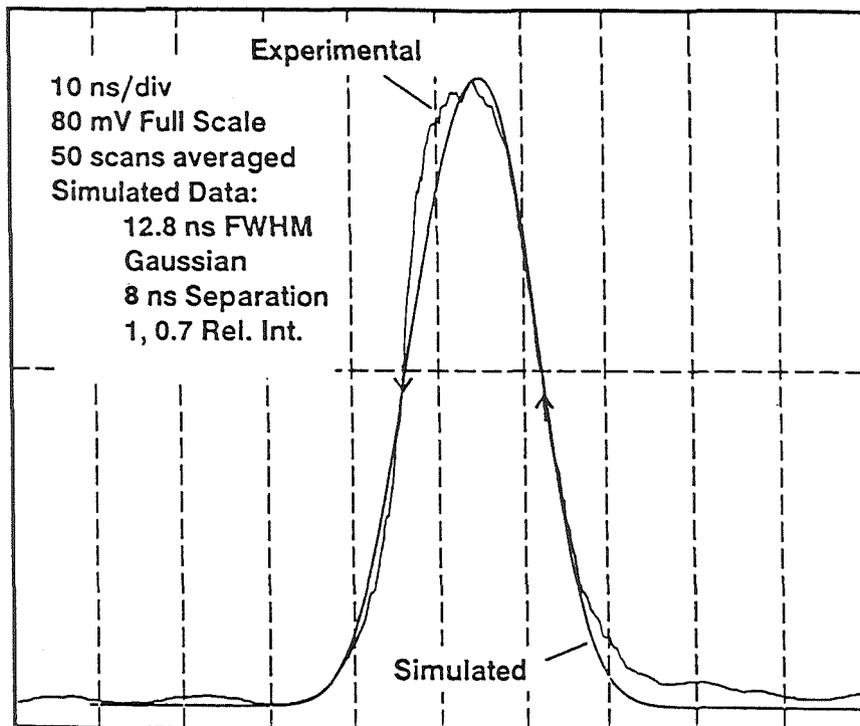


Figure 4. Experimental data (rough line) recorded with the carbon blocks spaced 1 foot apart. Simulated data (smooth line) consisting of two 12.8 ns FWHM pulses separated by 8 ns (*i.e.*, 4 feet due to *round trip travel*) with 1.0 and 0.7 relative intensities.

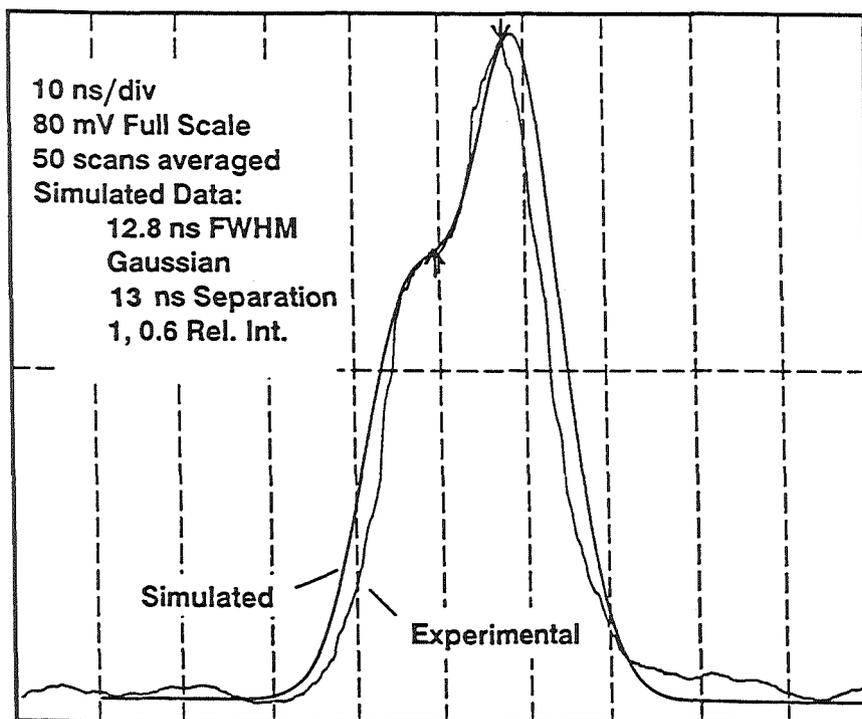


Figure 5. Experimental data (rough line) recorded with the carbon blocks spaced 3 feet apart. Simulated data (smooth line) consisting of two 12.8 ns FWHM pulses separated by 13 ns (*i.e.*, 6.5 feet *round trip travel*) with 1.0 and 0.6 relative intensities.

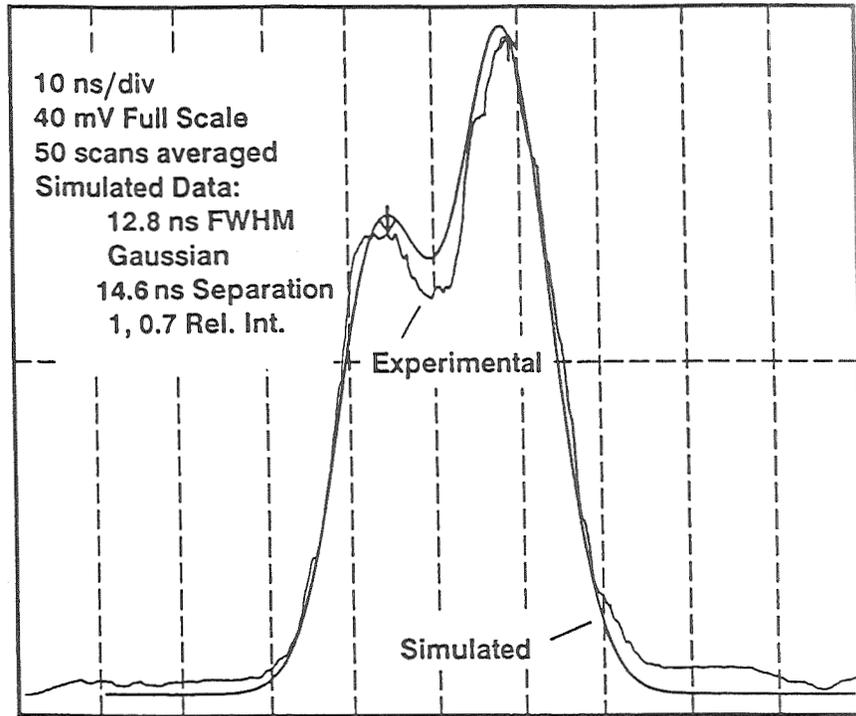


Figure 6. Experimental data (rough line) recorded with the carbon blocks spaced 5 feet apart. Simulated data (smooth line) consisting of two 12.8 ns FWHM pulses separated by 14.6 ns (*i.e.*, 7.3 feet round trip travel) with 1.0 and 0.7 relative intensities.

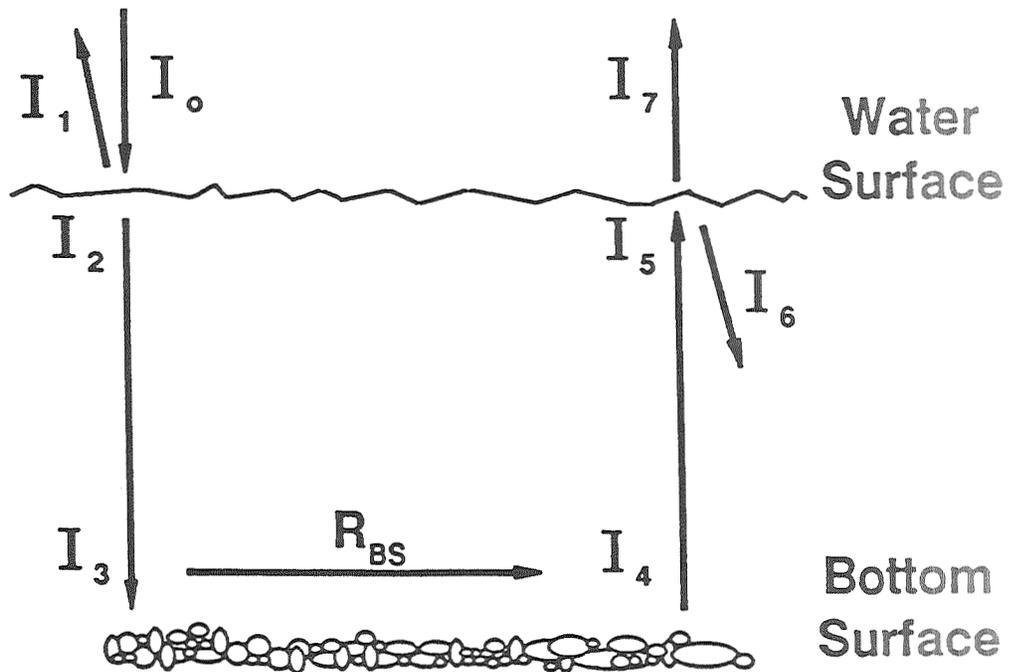


Figure 7. Simple model for determining relative return signal intensities for LIDAR bathymetry measurements.

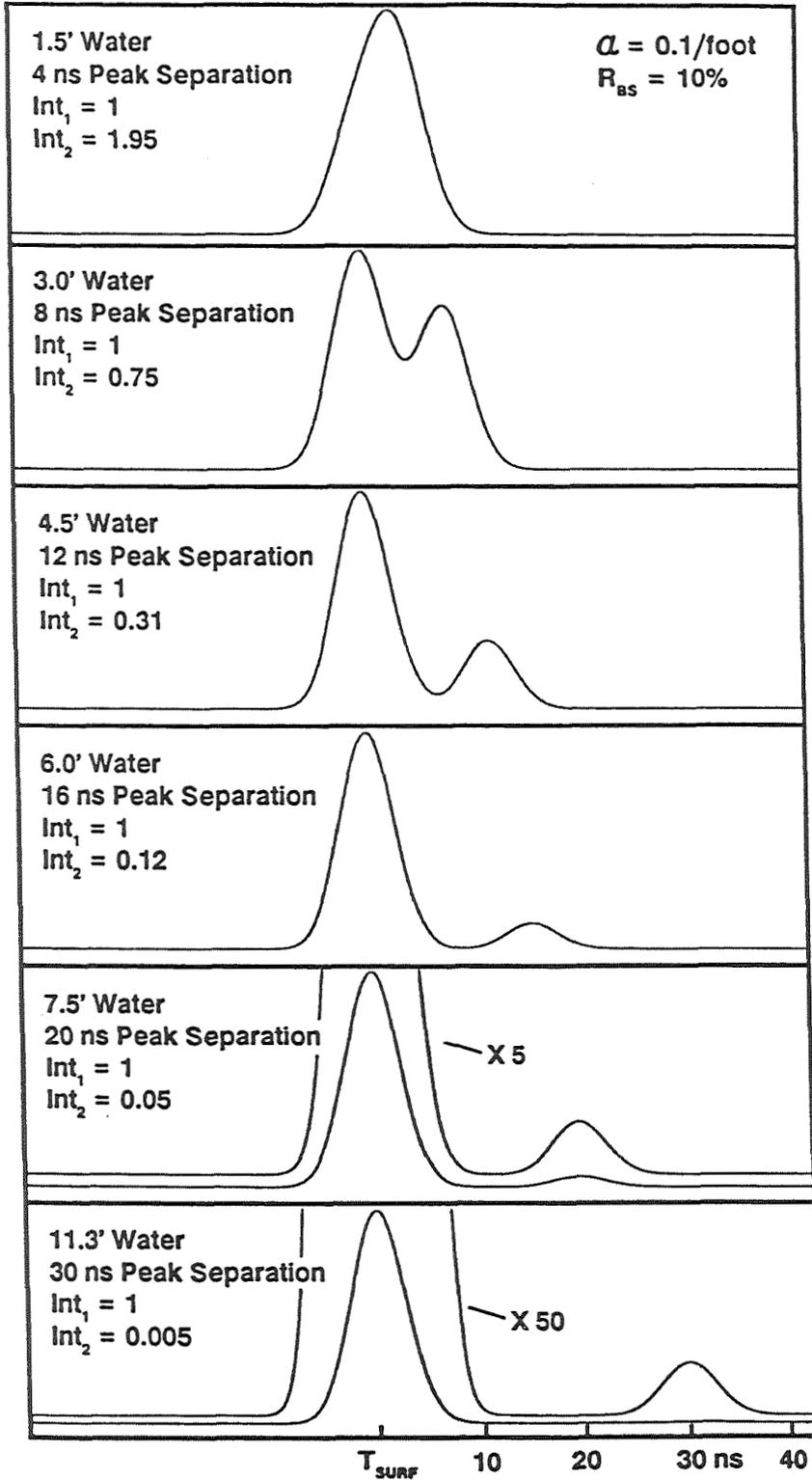


Figure 9. Family of simulated return waveform traces using a water attenuation coefficient (α) = 0.1/foot and bottom surface reflectivity (R_{BS}) = 10%. T_{surf} designates a reference time as the peak of the water surface return pulse.

COMMERCIAL APPLICATIONS MULTISPECTRAL SENSOR SYSTEM

N93-25615

Ronald J. Birk
Sverdrup Technology
Stennis Space Center, MS 39529

Bruce Spiering
NASA Science and Technology Laboratory
Stennis Space Center, MS 39529

55443
150524
p-12

ABSTRACT

NASA's Office of Commercial Programs is funding a multispectral sensor system to be used in the development of remote sensing applications. The Airborne Terrestrial Applications Sensor (ATLAS) is designed to provide versatility in acquiring spectral and spatial information. The ATLAS system will be a test bed for the development of specifications for airborne and spaceborne remote sensing instrumentation for dedicated applications. This objective requires spectral coverage from the visible through thermal infrared wavelengths, variable spatial resolution from 2-25 meters; high geometric and geo-location accuracy; on-board radiometric calibration; digital recording; and optimized performance for minimized cost, size, and weight. ATLAS is scheduled to be available in 3rd quarter 1992 for acquisition of data for applications such as environmental monitoring, facilities management, geographic information systems data base development, and mineral exploration.

1. BACKGROUND

The NASA Office of Commercial Programs (OCP) has a key cooperative effort with US industry to enhance our nation's economic standing in space-based remote sensing technology¹. OCP's Commercial Earth Observations Program (CEOP) is designed to increase US economic returns from remote sensing and related spatial information services. OCP provided Stennis Space Center (SSC) with mission requirements (Table I) for engineering and support of a versatile sensor system for developing remote sensing applications.

Table I Mission Requirements

(M1)	Provide prototype sensor system and subsystem test, verification, and applications proof-of-concept.
(M2)	Provide a flexible configuration for determining optimal sensor system design for specific applications to facilitate subsequent scanner development with private funds.
(M3)	Provide for calibration of other aircraft or satellite sensor systems during concurrent data acquisition missions.
(M4)	Utilize advanced and/or innovative technologies to decrease the cost and complexity of a remote sensing system, while attempting to increase the capability and end-user functionality of the data.

Stennis Space Center (SSC) acts as the lead center for remote sensing for OCP. SSC has operated an airborne multispectral data acquisition facility since 1972. The Thematic Mapper Simulator (TMS), Daedalus Enterprises^a

^a Mention of products/companies does not imply endorsement by the US Government, NASA, or Lockheed Engineering & Sciences Company. References are provided solely for the benefit of the reader.

developed Thermal Infrared Multispectral Scanner² (TIMS) and the Calibrated Airborne Multispectral Scanner (CAMS) have successfully supported 590 missions across the United States, Canada and Central America. These missions were predominately sponsored by NASA and other government agency earth resource monitoring projects. Although these sensors have some attributes to support commercial applications development, they are limited by the following:

- o Priority use to support science research projects.
- o System specifications designed to support science applications.
- o Configurations not conducive to spectral or spatial resolution modifications.
- o System/subsystem design based on > 5 year old technology.

Specifications for the ATLAS system were developed through an analysis of viable designs. These designs were based on functional requirements³ (Table II) derived from the mission requirements stated above and for compatibility with existing aircraft and on-going commercial programs at SSC.

Table II ATLAS Functional Requirements

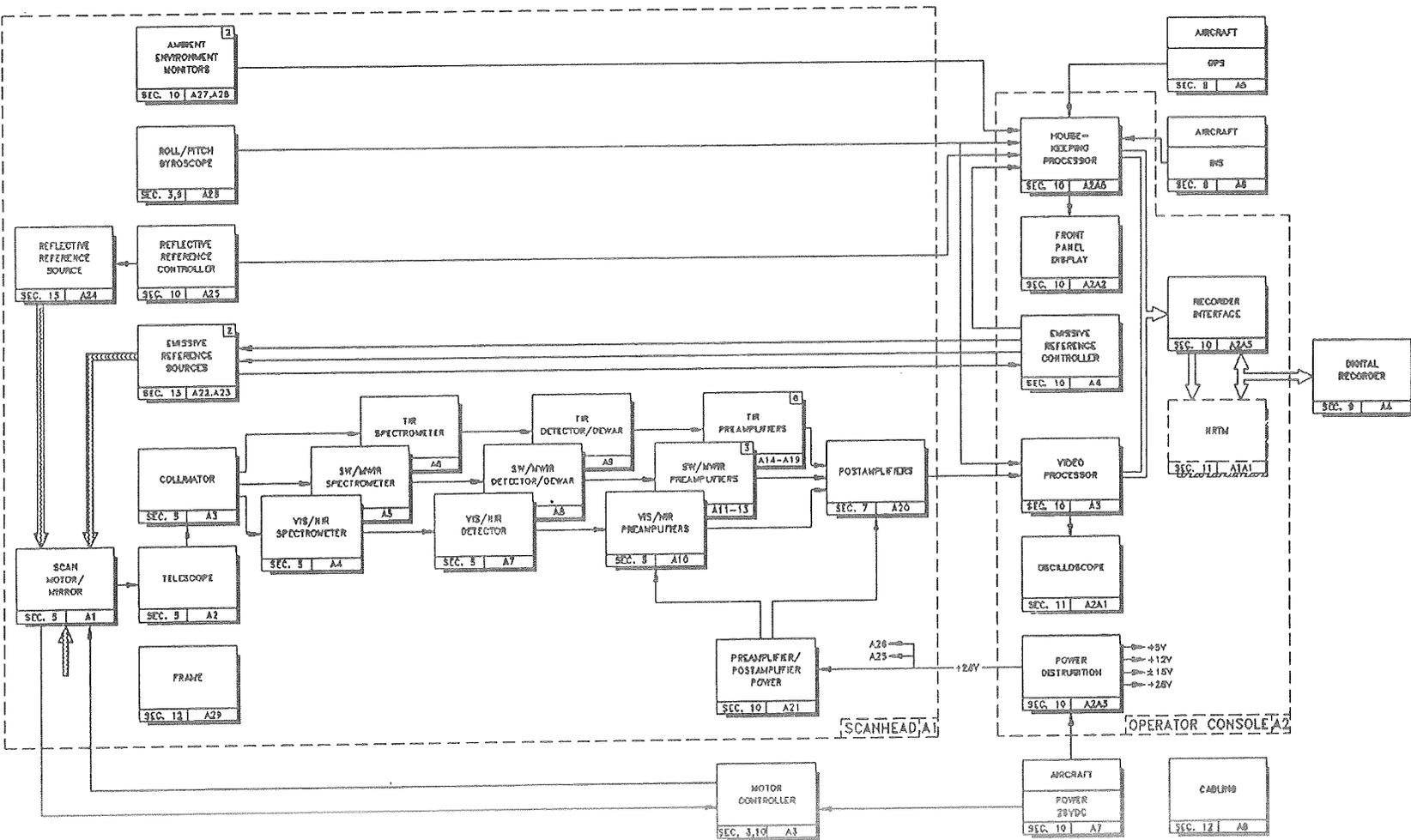
(F1)	Deployable from a Learjet 23 or compatible aircraft.
(F2)	Airborne and satellite-based sensor system emulation.
(F3)	Multispectral response from the visible through thermal infrared (0.4 - 12.5 μ m) wavelengths.
(F4)	Ground spatial resolution ranging from 2-25 meters, with design emphasis on less than 10 meters.
(F5)	Reflective and emissive calibration sources traceable to national standards.
(F6)	Modular and reconfigurable system subassemblies.
(F7)	Accurate aircraft attitude and geo-location knowledge to provide ancillary data for high geometric fidelity.
(F8)	Direct digital recording.
(F9)	High reliability, ease of maintainability, and operator interface optimization.

2. INTRODUCTION

This paper provides a functional description (Section 3) and performance parameters (Section 4) of the ATLAS system that is currently under development. Results of analysis⁴ to determine the system signal-to-noise ratio (SNR), noise equivalent radiance (NER) and noise equivalent temperature differential (NE Δ T) were calculated with a sensor analysis software package. The program is called Analytical Tools for Thermal Infrared Engineering (ATTIRE)⁵ and was developed at Stennis. Other system parameters include positioning accuracy, geometric fidelity, and digital recorder interface specifications. The technical specifications are listed in Appendix 6.2.

ATLAS (Figure 1) is an opto/mechanical scanner with multi-channel detector modules. ATLAS is designed to acquire data in 15 discrete channels in 5 atmospheric windows of the 0.4 - 12.5 μ m region of the electromagnetic spectrum. The system has 3 spectrometers referred to as: visible/near infrared (VNIR), shortwave/midwave infrared (SW/MWIR), and the thermal infrared (TIR). Simultaneous acquisition of data in is of significant value in evaluating optimal sensor system performance for a particular application. Channel center wavelength and

Figure 1 - Sensor System Block Diagram



bandwidth can be varied by changing modular spectrometers and/or detector assemblies.

Design of a visible through thermal infrared sensor system requires a detailed analysis of input signal propagation through the system. The major components in developing a model for the sensor system are the source, atmosphere, optics, detector, spatial parameters, and preamplifier electronics. The final goal of the analysis is to determine the performance parameters of SNR, NER, and NEAT for each channel. The atmospheric interference for each spectral bandpass is modeled with the PC version of LOWTRAN 7 and results are incorporated in the ATTIRE modeling of SNR, NER, and NEAT.

3. ATLAS SYSTEM

The ATLAS system consists of a scan head to be mounted in the equipment bay of the Learjet 23 aircraft and the operator console that is rack mounted in the cabin. The scan head (Figure 2) houses the scan motor/mirror assembly, telescope, 3 spectrometers, detector assemblies, calibration sources, and preamplifier electronics. The operator console includes the video processing electronics, ancillary data (housekeeping) control electronics, blackbody and motor controllers, and monitoring instrumentation.

3.1 Spatial Parameters

Spatial resolution is determined by aircraft altitude, velocity and scan speed. ATTIRE utilizes the relationship expressed in Equation 1 to determine scan speed effects on performance. An airborne platform capable of flying at altitudes between 1000 and 14,000 meters in altitude allows for sensor spatial ground resolution to be varied from 2-25 meters. The system has a variable integration time as a function of scan speed. To acquire 2 meter ground resolution, an aircraft velocity of 195 knots at an altitude of 1000 meters is required at a scan speed of 50 rps.

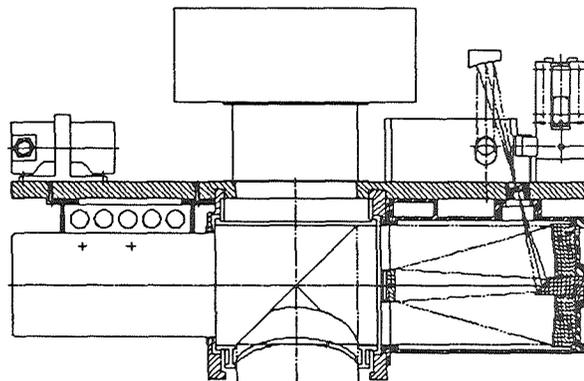


Figure 2 Scan Head Schematic

$$GIFOV = \alpha a = \frac{v}{n} \tag{1}$$

where	GIFOV	=	ground spatial resolution (meters)	n	=	scan speed
	a	=	altitude of aircraft (meters)	v	=	aircraft velocity (m/s)
	α	=	angular IFOV			

Distortions are introduced in the spatial data as a function of the scanning configuration (Figure 3). Examples include widening of pixels at the scan edges and the "S" effect due to the forward motion of the plane during a flight line and aircraft dynamics. These distortions may be corrected during post-processing with appropriate compensation data.

3.2 Optics

The optics subsystem collects, focuses, and disperses the radiant flux from the source to 3 spectrometers. The energy is then focused onto individual detectors. The two parameters that are critical to the collection of this energy are the area of the entrance aperture of the optics and the focal length. The radiance, limited by the solid angle subtended by the ground pixel, is incident on the entrance aperture and determines the irradiance available to the optics.

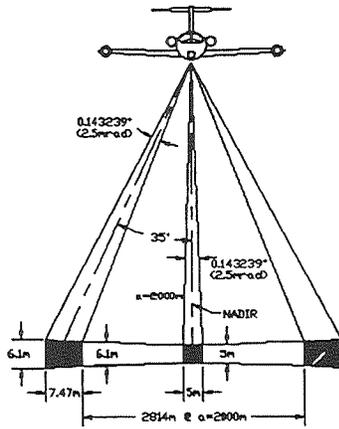


Figure 3 - Pixel Size Variation

25.5 mm and a back focal length (BFL) of 20.5 mm. The objective is well corrected over the 8-12 μm region and is ideally suited for this compact spectrometer.

The SW/MWIR spectrometer accepts the reflected radiation from the TIR dichroic filter and separates the 1.55-4.2 μm energy from the visible (0.45-0.90 μm) radiation using a second dichroic filter. The transmitted radiation is dispersed into 3 channels by a reflective grating onto an InSb detector array using an achromatic objective. The objective is a ZnS-Si-ZnS triplet that is 38 mm in diameter with an EFL of 32.5 mm and a BFL of 17.5 mm.

The VNIR spectrometer (Figure 4) disperses radiation between 0.45-0.90 μm into 6 channels. The VNIR radiation is separated from the longer wavelength energy by reflection from the second dichroic filter. A transmission grating with a groove spacing of 300 lines/mm is used to disperse the collimated radiation to a pair of achromatic doublets that focus the dispersed light onto the silicon detector array. The doublets are standard off-the-shelf achromats.

3.3 Detectors

The detectors were chosen based on performance in each spectral region. A 6 element silicon array will be used for the 6 VNIR channels. An indium/antimonide (InSb) array incorporated in a LN₂ dewar assembly will be used for the 3 SW/MWIR channels. A mercury cadmium telluride (HgCdTe) array incorporated in a LN₂ dewar assembly will be used for the 6 thermal channels.

3.4 Preamplifiers

The VNIR and SW/MWIR channels using photovoltaic (PV) detectors with a high impedance ($> 500 \text{ M}\Omega$) will use preamplifiers operated in an unbiased mode to take advantage of the characteristic high D^* value. The TIR channels using photoconductive (PC) detectors with low impedance (20-120 Ω) will have preamplifiers configured in a Wheatstone bridge circuit.

The gain requirements for each channel depend on the input signal from the detectors. The photovoltaic detector amplifiers will deliver a gain between 10^6 and 10^7 . The photoconductive detector amplifiers will have an average gain of 200-300. Bandwidth requirements for the amplifiers are established by the maximum scan speed of 50 rps. Therefore, the electronic bandwidth, Δf , must have a minimum value of 79.54 kHz.

The telescope and the 3 spectrometers were designed at SSC⁶ using Optical Research Associates (ORA) Code V optical design software. The telescope is a 7.5" aperture Dall-Kirkham design with an f-number of 7.84 and a 2 mrad field stop. Ground radiation is reflected by the scan mirror onto the telescope primary mirror. The converging radiation is reflected from the primary onto the secondary mirror, then to a folding flat mirror. The radiation is focused off-axis onto a square field stop aperture that defines the 2 mrad instantaneous field of view (IFOV) for the system. Above the field stop, an optical window is mounted to seal the spectrometer assembly from the atmosphere.

The TIR spectrometer uses a collimating and a folding mirror to direct the incoming radiation to a dichroic filter that separates the 8.2-12.2 μm radiation from the shorter wavelengths (0.45-4.2 μm). The transmitted radiation is dispersed into 6 channels by a reflective grating, and then focused onto a HgCdTe detector array using an achromatic objective. The objective is a germanium triplet of 40 mm diameter, with an effective focal length (EFL) of

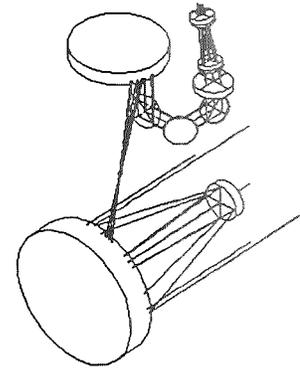


Figure 4 VNIR Spectrometer & Telescope

3.5 Video Processor

The major functions of the video processor are signal conditioning, anti-alias filtering, analog-to-digital conversion (ADC), and buffering. The analog signal from the postamplifiers on the scan head will be routed through fiber optic cable for EMI/RFI noise suppression to the operator console in the cabin. Gain and offset controls will be used by the operator to optimize dynamic range. Sufficient gain is applied to the signal to obtain a maximum amplitude of 10V p-p. The analysis to determine the digitization resolution resulted in a requirement of > 10 bits. To meet this requirement, a 12 bit Datel Inc. analog-to-digital convertor with integral sample and hold was chosen.

3.6 Timing

The position of the mirror is critical for precise timing of video sampling. The encoder resolution of the motor/mirror assembly is 4096 pulses per revolution, or 1.53 mrad. From this a master clock with a frequency of 12.5x the encoder frequency is generated. A timing diagram of the sampling sequence is shown in Figure 5.

3.7 Housekeeping

The housekeeping subsystem interface architecture (Figure 6) is designed to incorporate the ancillary data required during post-processing and provide other system functions such as roll compensation. Interfaces to the aircraft inertial navigation system (INS), global positioning system (GPS) receiver, blackbody controller, scan head gyroscope, ambient temperature sensors, front panel controls, real-time clock, mission specific data, and input for system calibration parameters are required. These functions are implemented through a dual microcontroller-based "embedded control" design. The microcontroller chosen was a Motorola 68HC11 running at a machine cycle of 0.5 μ s.

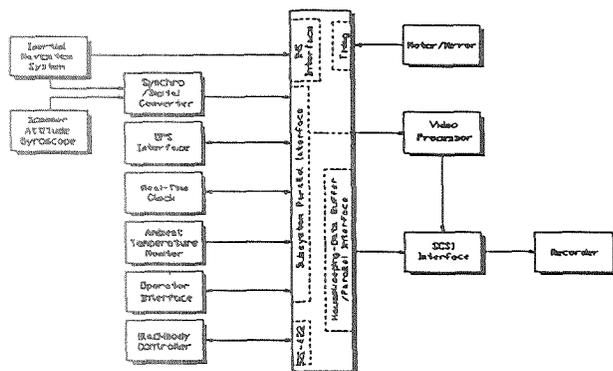


Figure 6 Subsystem Interface Architecture

The housekeeping subsystem interface architecture (Figure 6) is designed to incorporate the ancillary data required during post-processing and provide other system functions such as roll compensation. Interfaces to the aircraft inertial navigation system (INS), global positioning system (GPS) receiver, blackbody controller, scan head gyroscope, ambient temperature sensors, front panel controls, real-time clock, mission specific data, and input for system calibration parameters are required. These functions are implemented through a dual microcontroller-based "embedded control" design. The microcontroller chosen was a Motorola 68HC11 running at a machine cycle of 0.5 μ s.

A comparison of the accuracies of the aircraft inertial navigation system (INS) (values are taken at the equator) and differential GPS under best case conditions is presented in Table III. Two-dimensional (latitude and longitude) positional errors translate directly to the ground scene (assuming straight and level flight). A 5 meter lat/long

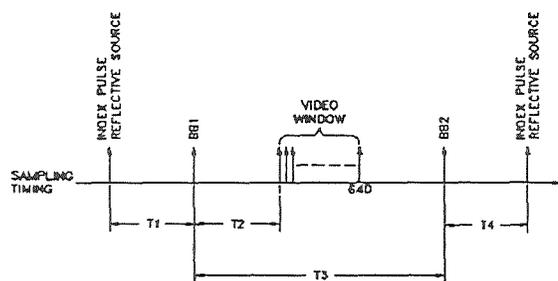


Figure 5 - Sampling Timing

3.7.1 Global Positioning System Interface

Geo-location knowledge is important for maximizing the usefulness of remotely sensed data. This is particularly true when the data are intended to be incorporated into a geographic information system (GIS). Accurate location knowledge is necessary for coregistering remotely sensed data sets to perform change detection analysis.

The ATLAS Global Positioning System (GPS) receiver will be a differential unit. Position, velocity, and time (PVT) is calculated at a 1-2Hz rate and is precisely time tagged to allow synchronization of the GPS signals with systems that are tied to it. The total number of

satellites in the GPS constellation will be 21 plus 3 active spares by 1994. The GPS constellation now consists of 13 satellites and provides world-wide 2-D coverage, with 3-D coverage in "windows" (times when four or more satellites are in view). An 18-satellite constellation will provide global 3-D positioning, and is expected sometime in 1992. The full constellation, in conjunction with differential operation, will allow for positional uncertainties of ≤ 5 meters in latitude, longitude, and altitude.

Table III GPS versus INS Accuracy

	LAT (m)	LONG (m)	ALT (m)	TIME (s)	DRIFT (°)	GND SP (knts)
INS	177	177	N/A	0.1	1	1
GPS	5	5	5	1e-9	N/A	0.3

circular error translates to 1 pixel for a 2 mrad system at an altitude of 4,125 feet. At an altitude of 41,250 feet, accuracies of 1/10 pixel are possible.

Geometric fidelity is a fundamental requirement of many applications of remotely sensed imagery. For airborne and spaceborne systems, geometric distortion is introduced in the imagery as a function of platform dynamics. An airborne scanner is subject to dynamic motion in roll, pitch, yaw, altitude, and velocity vectors which interact in a complex manner. Correction by manual intervention with standard processing techniques is time consuming and yields results with limited accuracy. Information provided by the GPS and system roll/pitch gyro provide for computer-assisted rectification of the data.

3.7.2 Roll/Pitch Gyroscope

A roll/pitch gyroscope is mounted directly on the scan head to minimize rotational error due to misalignment between the gyroscope and scan head. The outputs of the gyroscope are 3-wire synchro outputs that can be sampled continuously and are synchronized to the ATLAS system. The output of the aircraft INS gyro is also recorded.

The output is converted to a 14 bit digital number using two synchro-to-digital converters (SDC). Values for roll and pitch are stored in housekeeping. Roll output is used for real-time roll correction. The typical angular error for the gyroscope is $< \pm 0.2^\circ$ from true pitch and roll angles over full scale. This equates to ± 1.75 pixels at nadir for a 2.0-mrad system and represents the accuracy to which the gyroscope can locate the nadir pixel. It is a constant offset and once determined can be subtracted from all pitch and roll values. The angular repeatability is typically $\pm 0.05^\circ$ or less and is equivalent to ± 0.44 pixel. With the inclusion of quantization error, the angular error is within ± 1.93 pixels and the angular repeatability is within ± 0.92 pixel.

3.8 Digital Recorder

Recording of data using compact digital technology meets the mission requirement of implementing new technology to increase capability while decreasing cost. The 8 mm helical scan Exabyte digital recorder technology can handle most data rates generated by the ATLAS. The Exabyte 8500 recorder has a 5 Gbyte per tape capacity and a throughput of 0.5 Mbytes per second. Two units will handle 15 channels of 12 bit data at maximum scan speed. By designing the system to format the data to be compatible with image processing software and using a tape drive that complies to a standard computer interface (SCSI), data can be read directly into an image processing workstation upon delivery by the aircraft data acquisition crew. The 8 mm recorders are lighter (< 50 lb) and smaller (< 600 in³) than conventional wide-band analog units and have significantly lower procurement and operating costs. They eliminate the requirement for a decommutation system, which represents a savings of around \$1M in support hardware and software.

3.9 Operator Interface

The operator interface performs 4 basic functions:

- o Provide for operator control of system power, channel gain and offset adjustment, blackbody set points, scan speed, tape recorder, and oscilloscope.
- o Data entry of mission number, time, date and performance parameters.

- o Provide indicators of system operation and performance of dynamic range, noise, parameter values, image quality.
- o Conduct and monitor diagnostic tests.

The operator console is ergonomically designed for access to controls on the oscilloscope for monitoring analog, digital, and recorded signals. Provisions have been made to incorporate a near-real-time monitor (NRTM) to display video data in flight.

3.10 Calibration Sources

Radiometric calibration is required for performing quantitative data analysis. For a system that operates in both the reflective and emissive regions of the spectrum, calibrated reflective (lamp) and emissive (blackbody) reference sources, whose calibration is traceable to national standards, are required. The reflective source has been designed as a modified integrated sphere available from Labsphere, Inc. The blackbody sources were designed at SSC and manufactured by Mikron Industries.

4. SENSOR SYSTEM PERFORMANCE

ATLAS performance parameters (Table IV) were determined for each channel assuming a standard solar curve^{7,8} for the reflective channels, with an albedo of 0.3, and a 300°K blackbody curve for the emissive channels. The average values for atmospheric transmission from sea level to a nominal altitude of 7000 meters were determined from Lowtran 7 and tables from Hudson⁹. Optical transmission was determined for each channel through analysis in Code V optical design software and vendor transmission and reflection curves for individual optical components. The effective focal length varies for each channel; values are listed in Table IV. Detector D^* , was obtained from vendor test data and input as $3e13$ for the VNIR, $2.41e11$ for the SW/MWIR, and $4.46e10$ for the TIR. Preamplifier noise factor, n_r , determined in laboratory tests at SSC, was input as 1.2 and preamplifier bandwidth as set at 78.5 KHz. Scan speed was taken to be 25 rps, providing a dwell time of 12.73 μ seconds. The system performance parameters of SNR, NER, and NE Δ T were calculated using the ATTIRE system analysis modeling software. Results of the modelling of the system design using actual data that characterizes the subsystems indicates very good performance. NER for the visible channels is much less than 0.5 % and NE Δ T for the thermal channels is on the order of 0.2°C. A summary of the radiometric equations used to calculate the performance parameters is provided in Appendix 6.1. Technical specifications for the ATLAS system are presented in Appendix 6.2.

5. CONCLUSIONS

The ATLAS system will be accessible through Stennis Space Center for industry to use as a test bed for the investigation of potential remote sensing applications. Optimal spatial and spectral specifications can be determined through the acquisition of coincident spectral coverage over a range of ground resolutions. Specifications for a support a particular application can then be provided to commercial sensor manufacturers for production. This program is designed to facilitate the use of remote sensing technology to acquire information exceeding the current capabilities of the Landsat and SPOT sensor systems.

The sensor system will provide good image quality performance, high geometric fidelity, and direct sensor-to-computer interface through digital recording media. Commercial users have priority use of the ATLAS system to develop the use of remote sensing technology for applications such as environmental monitoring, facilities management, geographic information systems data base development, and mineral exploration.

Table IV ATLAS Performance Parameters

CH	Band Limit (μm)	Source ($\text{W cm}^{-2} \text{sr}_1$)	Atm (τ_a)	Optics (τ_o)	EFL (cm)	SNR (dB)	NER (%)	NE Δ T ($^{\circ}\text{C}$)
1	0.45-0.52	3.31e-4	0.60	0.17	61.0	92	0.0025	N/A
2	0.52-0.60	6.80e-4	0.60	0.30	64.5	98	0.0013	N/A
3	0.60-0.63	2.52e-4	0.60	0.30	39.4	93	0.0021	N/A
4	0.63-0.69	4.97e-4	0.60	0.30	56.0	96	0.0015	N/A
5	0.69-0.76	3.76e-4	0.60	0.24	60.5	93	0.0022	N/A
6	0.76-0.90	3.95e-4	0.60	0.20	89.8	90	0.0031	N/A
7	1.55-1.75	2.79e-4	0.80	0.35	36.0	35	0.5765	N/A
8	2.08-2.35	1.69e-4	0.80	0.32	40.4	30	1.0642	N/A
9	3.35-4.20	2.01e-3	0.80	0.13	66.4	47	0.458	0.118
10	8.20-8.60	3.99e-4	0.75	0.42	21.0	43	0.7327	0.19
11	8.60-9.00	4.69e-4	0.75	0.48	21.0	44	0.6234	0.16
12	9.00-9.40	5.07e-4	0.75	0.51	21.0	45	0.5775	0.15
13	9.60-10.2	7.04e-4	0.75	0.47	26.0	46	0.5151	0.14
14	10.2-11.2	9.51e-4	0.75	0.39	34.7	46	0.5089	0.13
15	11.2-12.2	5.29e-4	0.75	0.23	34.7	41	0.9142	0.24

6. APPENDICES

6.1 Radiometric Performance Analysis

This appendix presents calculations^{9,10} performed to determine SNR, NER, and NE Δ T for a channel of the ATLAS system .

6.1.1 Source Bandpass Flux

The bandpass flux exiting from the source pixel is calculated by

$$L_e(\lambda, T) = \int_{\lambda_1}^{\lambda_2} \frac{\epsilon(\lambda)}{\pi} \frac{c_1}{\lambda^5 (\exp(c_2/\lambda T) - 1)} d\lambda \quad \text{W cm}^{-2}\text{sr}^{-1} \quad (2)$$

where $\epsilon(\lambda)$ = spectral emissivity
 λ_1 = lower wavelength of the channel
 λ_2 = upper wavelength of the channel.

6.1.2 Bandpass Flux Incident on the Detector

This radiance propagates through the atmosphere and the optical path of the sensor system before it is incident on the detector. The radiance incident on the detector is thus attenuated by the atmospheric and optical transmittance, and can be calculated by

$$L_e^d(T) = \tau_a \tau_o L_e(T) \quad W \text{ cm}^{-2} \text{ sr}^{-1} \quad (3)$$

where τ_a = atmospheric transmittance
 τ_o = optical transmittance

6.1.3 Throughput/Power Incident on the Detector

The radiance incident on the detector multiplied by the throughput (Υ) of the system yields the power incident on the detector. The throughput is defined as the product of the area of the detector and the solid angle subtended by the detector towards the imaging lens. By the Invariance of Throughput Theorem, this throughput also equals the product of the pixel area and solid angle subtended to the collecting lens.

$$\Upsilon = A_{pix} \Omega_e = A_d \Omega_l \quad \text{cm}^2 \text{ sr} \quad (4)$$

where A_{pix} = Area of ground pixel
 A_d = Area of detector
 Ω_e = Solid angle subtended by ground pixel at entrance aperture
 Ω_l = Solid angle subtended by detector at the imaging lens
 imaging lens

6.1.4 Noise Equivalent Power (NEP)

Noise Equivalent Power is defined as the power incident on the detector such that the SNR is unity. Using D_{bb}^* , detector area, and electrical bandwidth, the NEP can be determined by

$$NEP(T) = \frac{\sqrt{A_d \Delta f}}{D_{bb}^*(T_B)} \quad W \quad (5)$$

where A_d = Area of the detector
 Δf = electronic bandwidth
 D_{bb}^* = D^* normalized to blackbody temperature

6.1.5 Signal-to-Noise Ratio (SNR)

An expression for the SNR of the signal from the detector may be obtained from the NEP calculation

$$SNR_d = \frac{P_d}{NEP(T)} \quad (6)$$

where P_d = Power to the detector

6.1.6 Noise Equivalent Radiance (NER)

NER is defined as radiance incident on the sensor that produces an SNR of unity. At the detector output, NER is

$$NER_d = \frac{L_c^{BB}(T)}{SNR_d} \quad W \text{ cm}^{-2} \text{ sr}^{-1} \quad (7)$$

6.1.7 Noise Equivalent Temperature Difference (NEΔT)

Noise Equivalent Temperature Difference is defined as the change in target temperature that produces an SNR of unity. The following is an expression for NEΔT

$$NE\Delta T = \frac{n_f}{\frac{P_d}{\sqrt{A_d \Delta f}} \frac{D^*(\lambda_{pk}) C_2}{\lambda_{pk} T^2}} \quad ^\circ K \quad (8)$$

where $C_2 = hc/k = 1.44e-4 \mu\text{m K}$

6.2 Specifications Summary

Table V ATLAS Technical Specifications

Optical		Geometric	
Entrance aperture	7.5"	GPS accuracy	5 meters
F-number (VNIR/SWIR/TIR)	2.46/1.2/0.95	INS accuracy	177 meters (at the equator)
Spatial		Gyroscope accuracy	0.05°
TFOV	73.34°	Roll correction	±15°
IFOV	2 mrad	Electronic	
Scan speed	6-50 rps	Digitization resolution	12 bits
Calibration Sources		Analog bandwidth	DC-78.54 kHz
Field Filling	100%	Video words	640/channel
Uniformity	>95%	Housekeeping words	200
Stability	>95%	Recorder	
Emissivity (thermal)	0.99	Throughput	0.8 Mbytes/s
Settability (thermal)	0.1°C	Storage Capacity	>2 Gbytes/tape
Accuracy (thermal)	0.1°C	Interface	SCSI
Weight	435#	Power	22-32 VDC @ 74 A

7. ACKNOWLEDGEMENTS

This study was performed for the Commercial Earth Observations Program at Stennis Space Center for the NASA Office of Commercial Programs. The authors wish to thank the engineering staff of the Advanced Sensor Development Laboratory for their outstanding work on the development of this sensor system and their contributions to this paper.

8. REFERENCES

1. Birk, R.J., Tompkins, J.M., and Burns, G.S. (1991) "Commercial remote sensing small satellite feasibility study," *SPIE Vol. 1495 Small Satellite Technology and Applications*, pp 2-12.
2. Palluconi, F., Meeks, G. (1985) "Thermal Infrared Multispectral Scanner (TIMS): An Investigator's Guide to TIMS Data," *JPL Publication 85-32*.
3. Birk, R.J., Christensen, E., and Alexander, T. (1991) "Airborne Instrument Testbed System (AITS) Strategic Plan," *Internal Report* (Contact R. Birk at Stennis Space Center for a copy).
4. ASDL Engineering (1991) "ATLAS System Analysis," *Internal Report*, (Contact Bruce Spiering at Stennis Space Center for a copy).
5. Jaggi, S. (1991) "ATTIRE (Analytical Tools for Thermal Infrared Engineering)," *Proceedings of the Second Annual JPL Airborne Geoscience Workshop*.
6. DaMommio, A. and Kuo, S. (1992) "Optical design for the ATLAS multispectral scanner," *SPIE Vol. 1690 Design of Optical Instruments* (Submitted for publication).
7. Thekaekara, M.P., Kruger, T., and R. Duncun (1962) "Solar Irradiance Measurements from a Research Aircraft," *Applied Optics*, Vol. 8 No. 8, pp 1713-1732.
8. Moon, P. (1940) "Proposed Standard Solar-Radiation Curves for Engineering Use," *J. Franklin Inst.*, Vol. 230, No. 5 pp 583-617.
9. Wyatt, C. (1987) *Radiometric System Design*, MacMillian Publishing Company, New York.
10. Hudson, R. (1969) *Infrared System Engineering*, John Wiley & Sons, New York.

INTERACTIVE FORECASTING
WITH THE
NATIONAL WEATHER SERVICE RIVER FORECAST SYSTEM

George F. Smith and Donna Page
Office of Hydrology
National Weather Service, NOAA
Silver Spring, MD 20910

355-47
150525
p-10

ABSTRACT

The National Weather Service River Forecast System (NWSRFS) consists of several major hydrometeorologic subcomponents to model the physics of the flow of water through the hydrologic cycle. The entire NWSRFS currently runs in both mainframe and minicomputer environments, using command oriented text input to control the system computations. As computationally powerful and graphically sophisticated scientific workstations became available, the National Weather Service (NWS) recognized that a graphically based, interactive environment would enhance the accuracy and timeliness of NWS river and flood forecasts. Consequently, the operational forecasting portion of the NWSRFS has been ported to run under a UNIX operating system, with X windows as the display environment on a system of networked scientific workstations. In addition, the NWSRFS Interactive Forecast Program was developed to provide a graphical user interface to allow the forecaster to control NWSRFS program flow and to make adjustments to forecasts as necessary. The potential market for water resources forecasting is immense and largely untapped. Any private company able to market the river forecasting technologies currently developed by the NWS Office of Hydrology could provide benefits to many information users and profit from providing these services.

INTRODUCTION

The U.S. National Oceanic and Atmospheric Administration (NOAA) is responsible for using science and service to manage the resources of the United States. The National Weather Service (NWS) supports this mission by providing river and flood forecasts and warnings for protection of life and property, and by providing basic hydrologic forecast information for environmental and economic well being. The Office of Hydrology (OH) supports NOAA's and NWS's missions through the design, development, testing, implementation, and support of a physically-based hydrologic forecasting system - the National Weather Service River Forecast System (NWSRFS).

In general, a river forecast system (or almost any system) can be viewed as having major components of (1) forces that drive the system, or data, (2) a mechanism to analyze the driving forces, or processing, (3) the heart of the system where the physical laws of motion are modelled, and (4) products of the system, or guidance information output for decision making. The relationships of these general functions of a river forecast system are shown in Figure 1. This paper will concentrate on the modelling and some output features which, as part of an ongoing OH project tied to NWS modernization, have been converted to an interactive, graphical form on computationally powerful scientific workstations. The paper will also describe how this technology could be used by the private sector to provide additional water resources forecasting services.

There are many components which together form the NWSRFS. The next section will present a brief background and history of the evolution of the NWSRFS, including some of the rationale for the existing structure which allows NOAA/NWS to have one of the premier river forecast systems in the world.

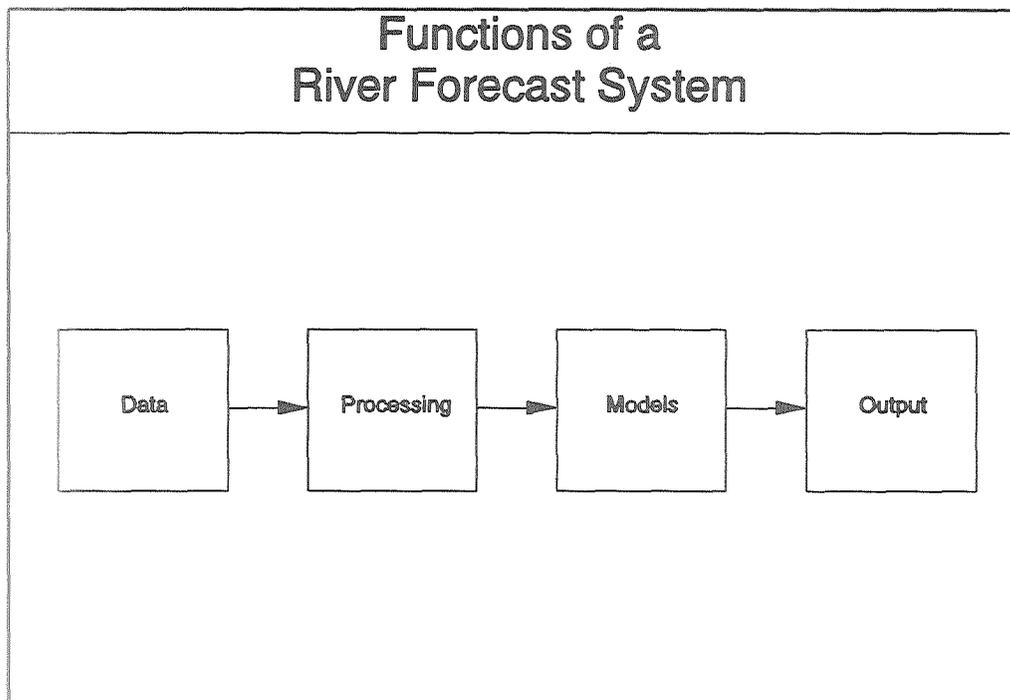


Figure 1

BACKGROUND/HISTORY

Prior to the advent and availability of digital computers many graphical or hand calculation methods were used for determining the flow of water in rivers. Because the hydrologic conditions varied greatly from one portion of the U.S. to another, different techniques for forecasting river conditions were developed by River Forecast Centers (RFC) responsible for different areas. There are presently 13 RFCs in the U.S. The areas of responsibility for the 12 which cover the coterminous U.S. are shown in Figure 2. The thirteenth RFC is responsible for the state of Alaska.

In the 1960's and early 1970's computers were introduced into the RFCs. Consistent with their pre-computer activities, each of the RFCs independently developed river forecasting software. Often this software was simply a computer representation of the graphical techniques used previously. These locally developed software programs introduced two major problems into the NWS forecasting activities. First, the forecasting software was dependent on the individual who did the initial development. When that person changed jobs or retired, much of the knowledge of how to run the programs, or how to maintain or enhance the programs was lost to the NWS. Second, forecasters at one RFC were trained in forecasting software that was, in general, only applicable to that RFC. If someone moved from one RFC to another they would have to be retrained in the forecast programs used at the new RFC. This also was a major burden to the NWS river forecasting mission.

In the early to mid 1970's the OH began development of the NWSRFS to (1) meet the forecasting needs of all RFCs, (2) be supported and documented at the National level, and (3) have enhancements and software configuration management coordinated by OH. One of the initial goals was to design a system which included existing techniques from many of the RFCs so that a single system could be used for river forecasting throughout the U.S.

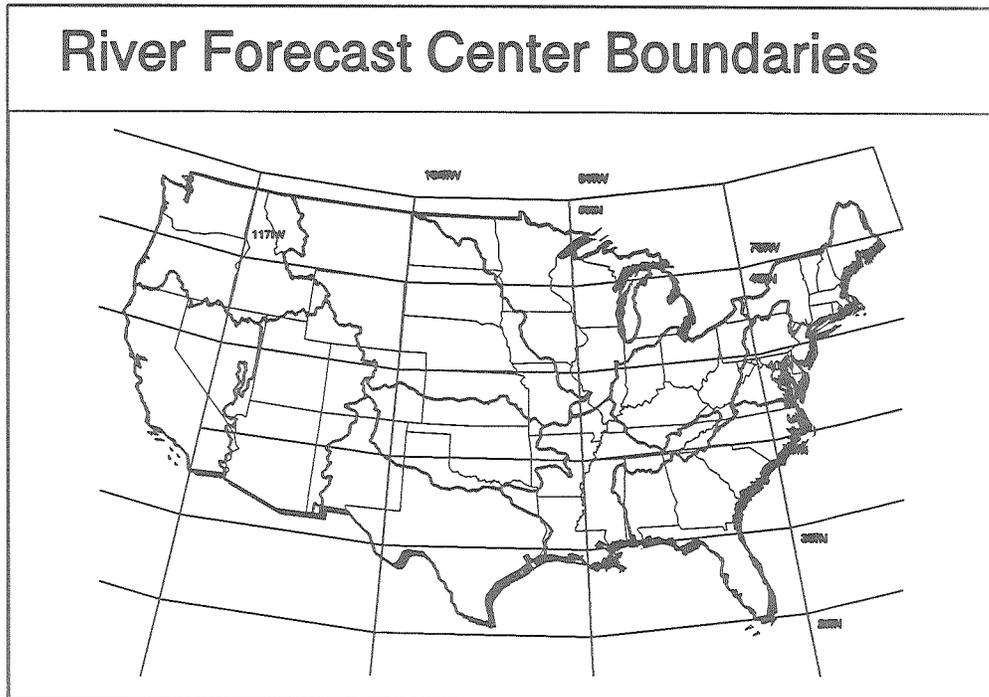


Figure 2

In the mid to late 1970's, initial versions of the NWSRFS were developed by software contractors under guidance from OH. These initial versions met some of the intended requirements of a national river forecast system, but they suffered from several basic flaws. Early versions of NWSRFS did not include all the features needed to model the flow in rivers in the varied hydrometeorologic regimes found throughout the U.S. Also, they did not account for the growth and evolution of computer technology and advances in hydrologic science. Versions 1 through 4 of the NWSRFS had a rigid program structure which made it difficult to add new modules as additional features were developed. The hydrologic modelling structure required that all basins use the same models in a fixed sequence. With the hydroclimatic variation found in the U.S. from humid to arid, and snow to sub-tropical conditions, this restriction was very limiting. New models or technology were very difficult to add to these early versions of the NWSRFS.

NWSRFS VERSION 5

In 1979, the OH began a project to completely redesign the NWSRFS. In addition to fixing the shortcomings found in previous versions, a major objective of the project was to develop a system structure which looked toward the future of hydrometeorologic forecasting. The initial requirements for NWSRFS Version 5 were developed from extensive interactions between designers in OH and the RFCs. Version 5 differed from previous ones in several ways, a major one being that scientific algorithms were designed to be independent of any specific computer system, and were coded by OH and RFC hydrologists who were intimately familiar with the physics of the processes being modelled. Specifications for data access and command decoding routines were developed by OH and RFC staff, and were coded by software contractors. The functional requirements which guided the design of NWSRFS Version 5 were to:

1. allow for a variety of models and procedures,
2. let the user control selection of models and sequence of use,
3. easily add new models and procedures to keep up with technological changes,
4. efficiently process large amounts of data to produce forecasts at hundreds of locations for each RFC, and
5. allow the user to flexibly control real-time processing.

Version 5 was designed to be modular, so that components could be developed by a number of individuals and then combined into a total system. References in the program code to system specific routines were isolated so that the entire NWSRFS could be ported from one hardware/operating system platform to another with minimum effort. Routines which performed scientific algorithms were separated from input/output routines so that the science could be run on any computer without needing changes in the reading or writing of information from the computer system. Scientific algorithms were organized into modular functions so that the functions could be shared, unchanged, among major components of the NWSRFS.

The functions representing one scientific algorithm, such as a snow, soil moisture, or river routing procedure are called an operation. In general, an operation in the NWSRFS is a set of functions that performs actions on a time series. Typically an operation describes the equations of motion governing the flow of water through a portion of the hydrologic cycle. There are also operations to display results, or to perform utility functions such as adding two time series. Table 1 provides a list of some of the currently available operations in the NWSRFS.

Table 1. NWSRFS Hydrologic Models

Snow	HYDRO-17 Snow Model
Soil	Sacramento Soil Moisture Accounting Ohio RFC API Rainfall-Runoff Model Middle Atlantic RFC API Rainfall-Runoff Model Central Region RFC API Rainfall-Runoff Model Colorado RFC API Rainfall-Runoff Model Xinanjiang Soil Moisture Accounting Continuous API Model Middle Atlantic RFC API Rainfall-Runoff Model #2
Channel	Channel Loss Dynamic Wave Routing Lag and K Routing Layered Coefficient Routing Muskingum Routing Tatum Routing Stage-Discharge Conversion Single Reservoir Simulation Model Unit Hydrograph

The operations that model the flow of water through the hydrologic cycle fall generally into the categories of (1)

one location to another on a river. Operations form the scientific heart of the NWSRFS and are shown in Figure 3 to be shared by the major sub-systems which comprise the NWSRFS Version 5. Because of the modular nature of the functions which make up any operation, functions can be shared with no change whatsoever among the programs which form the NWSRFS. This also allows new scientific techniques to be developed in the structure specified for an operation, and once tested to be immediately available for use in forecasting with the NWSRFS.

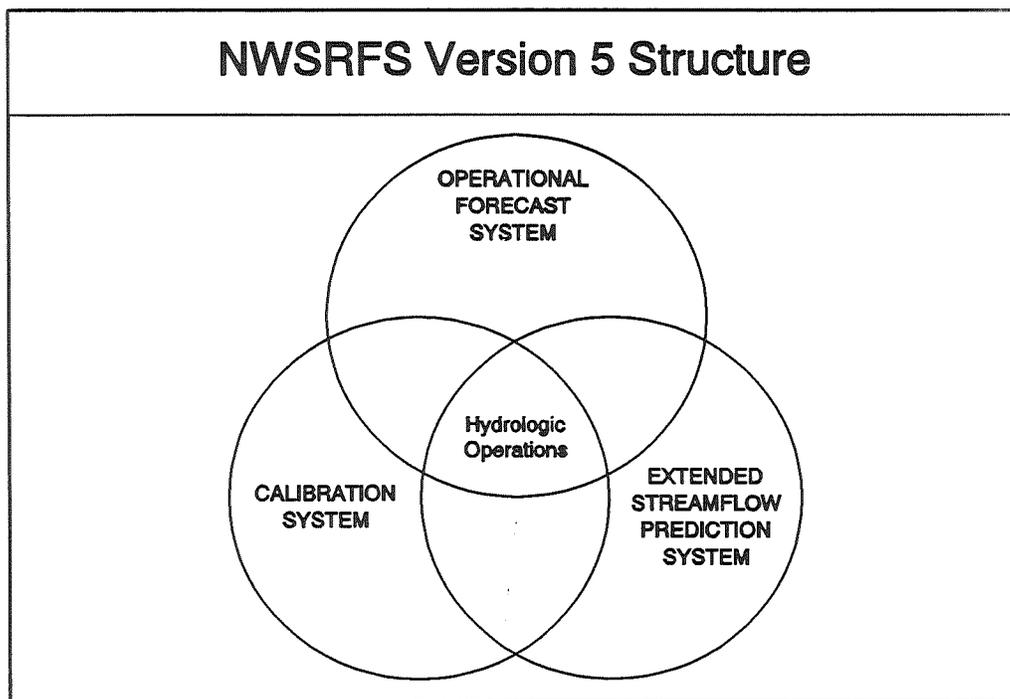


Figure 3

Hydrologic operations in NWSRFS are organized into an "operations table" to specify the physics of water movement for any subbasin. Operations can be selected from the list shown in Table 1. The order in which they are computed depends on the hydrometeorologic conditions of the subbasin being modelled. RFC forecasters can use their hydrologic expertise to determine the best sequence of scientific algorithms (operations) to model each subbasin. In this way, NWSRFS provides a generalized river forecasting system which can be used to model basins in any hydroclimatic regime. An example of the specific operations table for the Tahlequah, Oklahoma subbasin in the Arkansas-Red Basin RFC area is shown in Figure 4.

Initial NWSRFS Version 5 development occurred from 1979 through 1984. In 1985 NWSRFS Version 5 was delivered to the Arkansas-Red Basin RFC for initial operational forecasting use. Since then Version 5 has been installed in other RFCs and has been used daily to produce operational forecasts at thousands of locations along rivers throughout the U.S. New subbasins are continuously being calibrated and added as operational forecast locations by RFC hydrologists. Many new scientific algorithms and enhancements to existing operations have been added to improve the hydrologic modelling capabilities of the NWSRFS.

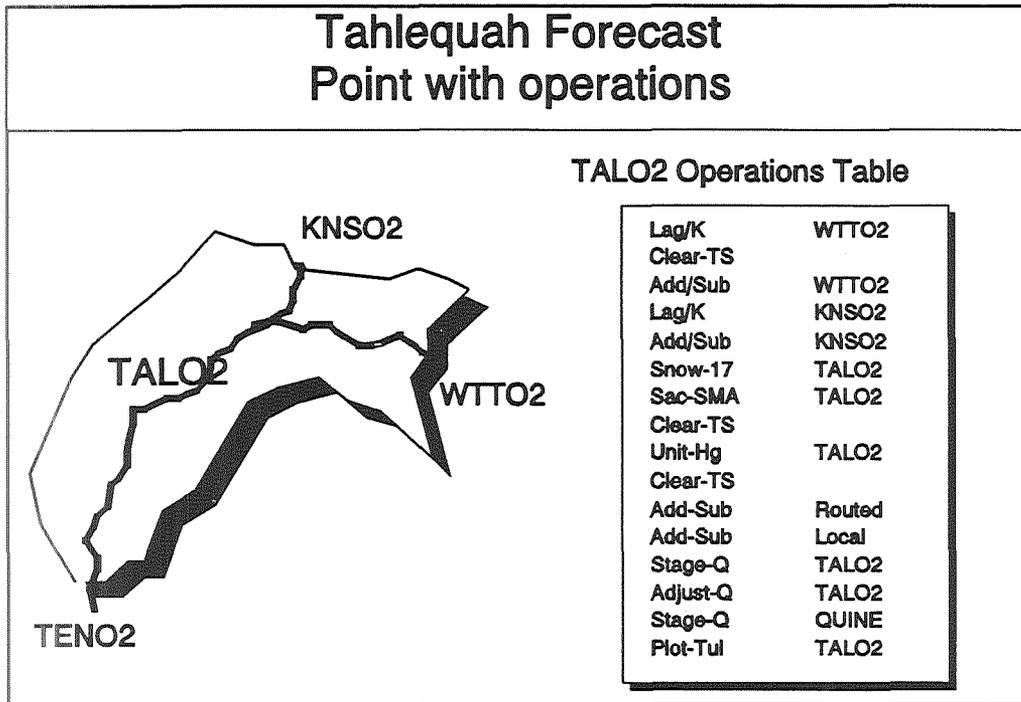


Figure 4

As computer technology has evolved the NWSRFS has kept pace. The initial NWSRFS design and development was on mainframe computers (NAS 9000s) at the NOAA Central Computer Facility (CCF). As minicomputers became powerful enough to support the system requirements of the NWSRFS, the NWSRFS Operational Forecast System (OFS) was ported to Prime minicomputers which are at OH and several of the RFCs. With the explosive growth in computational capabilities for scientific workstations, OH initiated a project in the late 1980's to prepare for modernization of the entire NWS by moving the scientific operations and forecasting component of the NWSRFS onto IBM RS/6000 workstations.

When the NWSRFS is run from the NOAA CCF, command input is sent over Remote Job Entry (RJE) lines from RFCs to the CCF. Line printer results are sent back to the RFC for display on standard printers or on text display screens.

Beginning in 1989, graphical display and user interface capabilities were developed for the NWSRFS. The result is the NWSRFS Interactive Forecast Program (IFP) which will be discussed in more detail in the next section of this paper.

INTERACTIVE FORECAST PROGRAM

The process of hydrologic forecasting requires human-machine interaction. This is because:

1. the equations with which we represent the physics of the hydrologic cycle do not perfectly

2. model the actual movement of water,
2. the process we use to calibrate, or find specific parametric values for, the models does not produce perfect results, and
3. we do not perfectly observe rainfall or stream conditions as input to the models.

In order to properly forecast a hydrologically connected series of subbasins, a forecaster must make decisions for each location along the river where observed river conditions are available. If values simulated by NWSRFS do not agree with observations, the forecaster must decide on the most likely source(s) of error, and make adjustments. When a river system is forecast with NWSRFS on the NOAA CCF or a Prime minicomputer, a group of subbasins are processed in a single batch run. Errors in upstream subbasins propagate into downstream basins, making forecasts for those basins less reliable. The only way to avoid this problem is by making adjustments to reduce or remove the error in any subbasin before processing downstream subbasins. The NWSRFS IFP provides the forecaster with this capability. An additional benefit of the IFP is the enhanced display capabilities of high-resolution color display terminals above those of line printer output.

As described above, hydrologic forecasting is inherently interactive. The initial designers of NWSRFS recognized this, but were limited because computational requirements demanded that the forecast system run on a mainframe computer with little interactive capabilities. The computational capabilities of scientific workstations have evolved so that the initial design features of NWSRFS Version 5 to allow for interactive forecasting can be realized.

Graphical user interface (GUI) and graphical display capabilities were developed on scientific workstations. Figure 5 shows in heavy outlines those portions of the mainframe and minicomputer versions of NWSRFS that were ported to scientific workstations and linked with the GUI and graphical display modules. The division of components among those solely in the NWSRFS OFS, those solely in the IFP, and those shared by both programs is shown in Figure 6.

Important features of the NWSRFS IFP include:

1. an operationally proven set of hydrologic models,
2. a system configuration which uses the UNIX operating system with X Windows graphical display protocol and Open Software Foundation (OSF) Motif,
3. adherence to OSF standards to be computer hardware platform independent,
4. a GUI that provides easy, powerful user interactions,
5. scientific applications that are isolated from the operating system specific function calls and input/output, and
6. the use of both C and FORTRAN programming languages; C for user interface and graphical display routines, FORTRAN for physical process modelling.

The IFP currently runs in two configurations, depending on the equipment available at a site. In the first, a Prime minicomputer runs the NWSRFS OFS and creates a current set of model conditions and time series. A forecaster at a scientific workstation networked to the minicomputer begins an IFP session by asking for information about a set of subbasins. This initial information is transferred from the minicomputer to the workstations. The remainder of the IFP session is performed on the workstation with computations of the operations tables for subbasins being forecast, adjustments made through the IFP GUI, and results displayed for forecaster interpretation. At the end of an IFP session, adjustments made for any subbasins are transferred to the minicomputer to become incorporated in further forecasting activities.

In the second configuration (Figure 7), a UNIX based fileserver replaces the Prime minicomputer. This eliminates the need to transfer information between different operating system environments and allows the NWSRFS OFS and IFP to operate more efficiently.

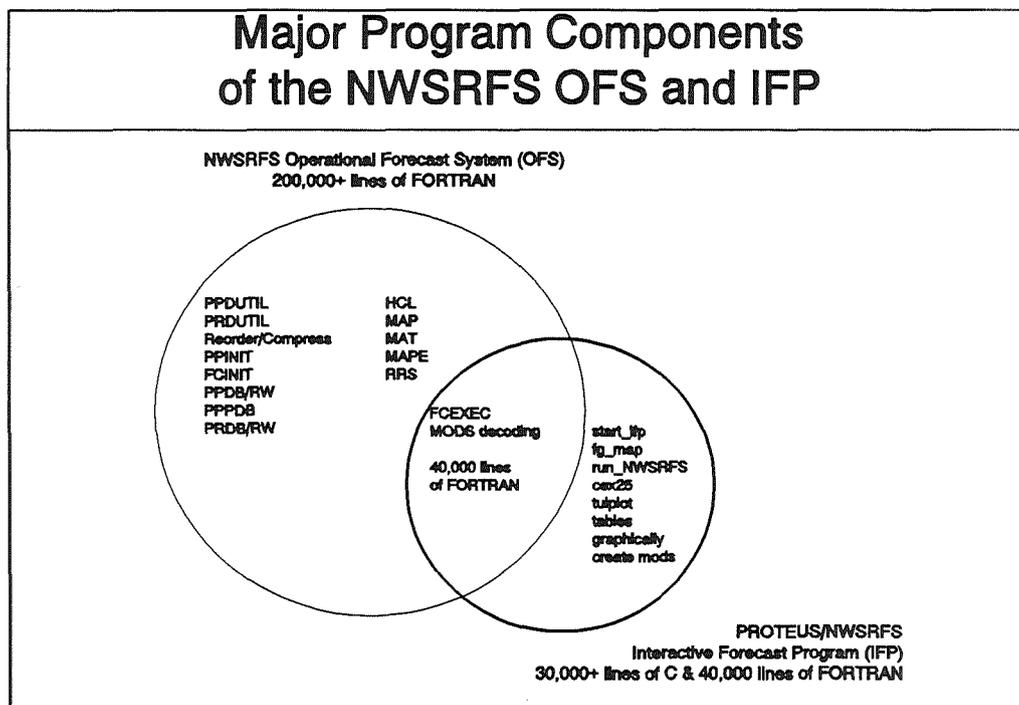


Figure 6

COMMERCIAL POTENTIAL

The demand for hydrologic forecasting services has been demonstrated by the response of the private sector to two subcomponents of the NWSRFS that are currently marketed by software engineering firms. The software firms have been successful in their sales of the dynamic wave river model and the dam break flood prediction model they have repackaged.

The potential market for timely and accurate hydrologic forecasts is immense and largely untapped. Navigation and recreation interests need information to help efficiently determine appropriate activities and resources for the river systems on which they operate. Reservoir operators have a tremendous potential gain by using river flow forecasts to help balance the competing needs of water supply, flood control, hydroelectric power generation, and ecological viability of the river.

An example of how these NWSRFS OFS and IFP technologies can be used commercially is provided by the Bonneville Power Authority (BPA) which operates a series of reservoirs on the Columbia and Willamette Rivers in the northwestern United States. Among the goals of their reservoir operations are flood mitigation, power generation, and maintenance of the river ecology to support the fishing industry in the northwest. In order to optimize these interests, BPA decision makers must, at times, balance the survivability of fish fingerlings with profits from power generation in their reservoir operations.

The information provided by the NWSRFS OFS and IFP would be useful to BPA in making those decisions.

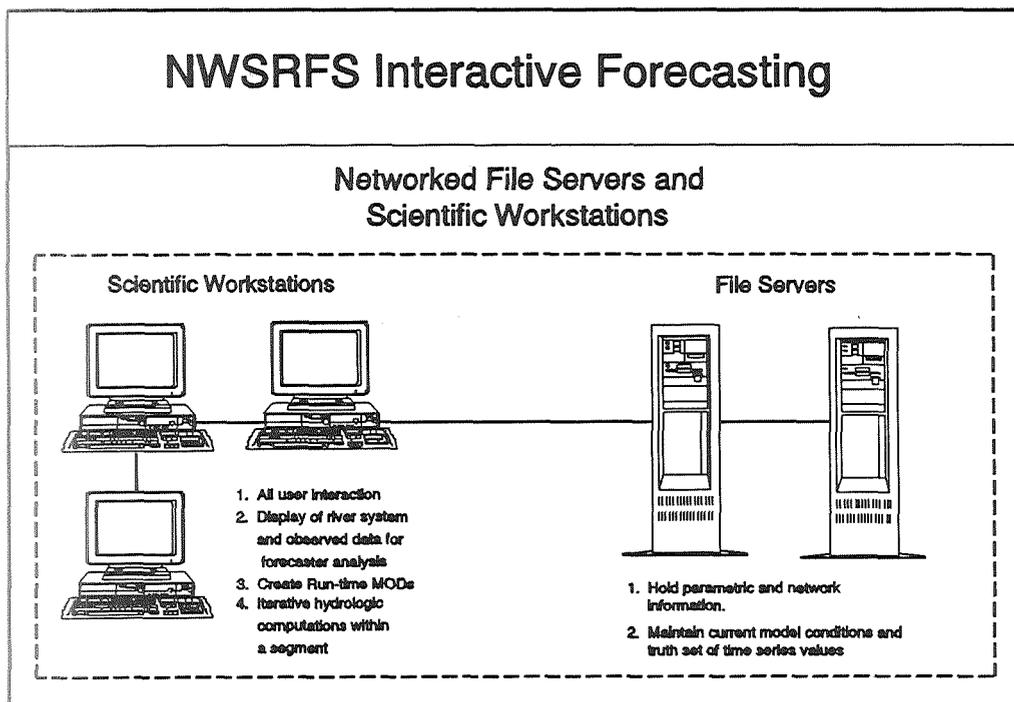


Figure 7

The IFP provides a graphical display of flows into the reservoirs and along the river system. These flows can be used to estimate power generating capabilities and track favorable fish migration conditions. The IFP also provides a flexible, easy-to-use interface to the powerful capabilities of the NWSRFS OFS that allows forecasters to try numerous what-if scenarios to visualize the effects of their reservoir operations.

The BPA is currently working with a private engineering firm to provide this information to their decision makers to optimize revenue from both fishing and power generating activities. This is just one example of the commercial potential of the NWSRFS OFS and IFP technologies. Any private company able to market the river forecasting technologies currently developed by the NWS Office of Hydrology could benefit any of the information users mentioned above (and many others) and profit from providing these services.

**SPACE LIFE SUPPORT TECHNOLOGY APPLICATIONS TO
TERRESTRIAL ENVIRONMENTAL PROBLEMS**

Steven H. Schwartzkopf
Howard L. Sleeper
Lockheed Missiles & Space Co., Inc.
Sunnyvale, CA 94088

566-54

156526

p. 7

ABSTRACT

Many of the problems now facing the human race on Earth are, in fact, life support issues. Decline of air Quality as a result of industrial and automotive emissions, pollution of ground water by organic pesticides or solvents, and the disposal of solid wastes are all examples of environmental problems that we must solve to sustain human life. The technologies currently under development to solve the problems of supporting human life for advanced space missions are extraordinarily synergistic with these environmental problems. The development of these technologies (including both physicochemical and bioregenerative types) is increasingly focused on closing the life support loop by removing and recycling contaminants and wastes to produce the materials necessary to sustain human life. By so doing, this technology development effort also focuses automatically on reducing resupply logistics requirements and increasing crew safety through increased self-sufficiency. This paper describes several technologies that have been developed to support human life in space and illustrates the applicability of the technologies to environmental problems including environmental remediation and pollution prevention.

INTRODUCTION

On previous missions, spacecraft life support systems have used open-loop, non-recycling chemical and mechanical technologies. These technologies were simple and sufficiently reliable to support humans for missions of relatively short duration. For NASA's planned Space Exploration Initiative (SEI), however, life support technologies must address a new and different set of requirements. These include longer mission durations and the need for closure of the life support loop (through material recycling) to minimize resupply logistics and maximize crew safety. Meeting these requirements involves new methods of life support which emphasize regenerative technologies.

The development of these regenerative life support technologies holds a promise which extends far beyond the SEI program, however. Many of the problems the human race must grapple with on Earth are fundamentally life support issues. Atmospheric changes caused by air pollution, and which contribute to the greenhouse effect or to depletion of the ozone layer present critical challenges to the agricultural and ecological sciences. Pollution of water by fertilizers, organic pesticides or chemical solvent presents serious health problems to our population. The recycling of solid wastes is another example of a problem for which we must find an appropriate technological solution. The objective of this paper is to describe some technologies currently under development for life support applications, and to outline how they might be applied to help solve some of these pressing environmental problems.

LIFE SUPPORT FUNCTIONS

Human beings require substantial amounts of material to sustain life. Including the water required for showers, personal hygiene, and food preparation, without recycling it takes over 3 and one-half tons of food, water and oxygen to support an average person for one year. If clothes wash water is added, the total weight of required life support materials more than doubles. Clearly, recycling life support materials is essential to minimize resupply for space missions.

The basic functions necessary to support human life include water recycling, solid waste processing, atmosphere regeneration, and food production. This paper presents examples of environmental applications

of technologies which fulfill the first two of these functions. The technologies available to provide basic life-support functions fall into three categories: physicochemical, bioregenerative, and hybrid. Physicochemical technologies are those which use chemical reactors and/or mechanical devices (such as fans, pumps, and filters). Bioregenerative technologies incorporate living organisms such as plants, algae, or bacteria to supply specific life support functions. Hybrid systems are created by combining components of both the technologies that are best suited to perform specific life support functions. Because of the clear need to reuse materials within the life support system, we are also focused on regenerative technologies which support recycling.

WATER RECYCLING

The recycling of water is one of the most significant challenges facing life support designers, because of the prodigious amounts of water used by people. Two types of water require processing for recycling, gray and black. Gray water is the effluent resulting from basic hygiene activities, food preparation, and collection of condensate from the atmosphere. Black water carries urine and fecal waste materials.

The physicochemical technologies developed for recycling water include simple distillation, filtration (e.g., reverse osmosis, hyperfiltration) and phase change processes (e.g., Vapor Compression Distillation (VCD), air evaporation). The bioregenerative processes used most frequently in water recycling are bacterial filters. Water reclamation by communities of higher aquatic plants and their associated bacteria is a relatively new method of bioregenerative processing (Wolverton, et. al., 1983). Waste water is pumped through a bed of aquatic plants which, along with the bacteria around their roots, remove contaminants. This type of system is being evaluated for tertiary processing of sewage water in several cities (e.g., San Diego).

Commercial waste water treatment plants generally use a hybrid approach incorporating both physical and bacterial filtration to treat water before discharging it into the environment. Such treatment does not normally remove all of the contaminants, and as a consequence the treated water is not purified sufficiently well to recycle directly for drinking, washing, or cooking.

Lockheed, in conjunction with Louisiana State University, is currently conducting research and advanced development work in which selected bacterial species are used to purify water for direct recycling (Miller, et. al., 1992). By matching the contaminant removal capabilities of bacterial species to the contaminants in the water, high levels of water purification can be achieved. Several such bacterial systems are now in operation, and some are being applied in site restoration efforts to purify polluted ground water supplies. In such applications, ground water is pumped through the bacterial reactor. After the bacteria metabolize the contaminants, the cleaned water may be returned to the aquifer.

Figure 1 provides a block diagram of an immobilized bioreactor for studying bacterial degradation of organic contaminants in the laboratory. Figure 2 illustrates both the laboratory and field (site remediation) bioreactors. The laboratory reactor is configured to support degradation studies in either plug flow or recycle configurations. Figures 3 and 4 show typical performance data for this type of laboratory reactor. In Figure 3, the degradation of two different concentrations of phenol in water was evaluated. As this figure shows, removal at the high phenol concentration was extremely high, averaging over 99.5% during the test period. Removal at the low phenol concentration averaged about 97% over the 8-day test.

Figure 4 shows degradation of three chlorinated hydrocarbon species during long duration tests of the bioreactor. Minimum efficiency was about 60% for TCE a 12ppm in water, but average efficiencies over the 25 day period were in the 80-95% range. Operation over long period of time may require addition of nutrient supplements to ensure the bacteria maintain a high operating efficiency.

WASTE PROCESSING

Due to the short durations of previous missions, waste processing has had little need to advance beyond the technologies for collection and storage. On the Mercury through Apollo missions, human wastes were stored

in holding compartments. Germicides were added to inhibit bacterial degradation of the wastes. Both Skylab and Shuttle dry and store organic waste materials. Other solid wastes (food and drink packages, tissues and wipes, etc.) are normally bagged and returned to Earth.

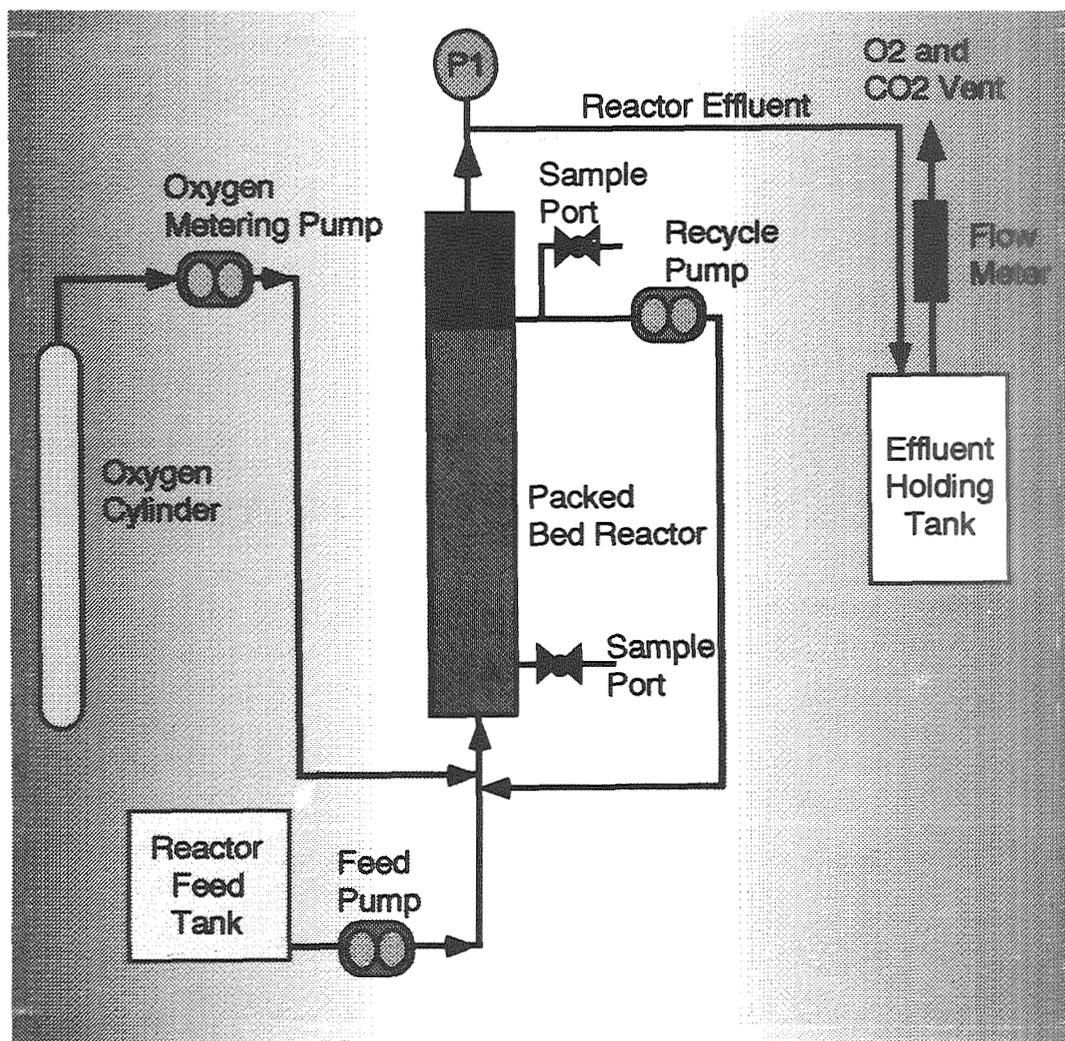


Figure 1. Block diagram of immobilized bioreactor

Physicochemical technologies for waste processing include incineration, electrochemical oxidation, wet oxidation, and supercritical wet oxidation. These systems typically require large amounts of electrical energy, but many highly toxic compounds can be effectively broken down by the intense physical conditions they produce.

Bioregenerative technologies include bacterial reactors (both aerobic and anaerobic) and combination higher plant/bacterial systems. Aerobic bacterial systems typically require higher energy inputs to maintain oxygenation (e.g., aerating pumps, mixers). Anaerobic systems require very little energy, but have very slow process rates, and the anaerobic bacteria are more susceptible to changes in environmental conditions (Wolverton, et. al., 1983).

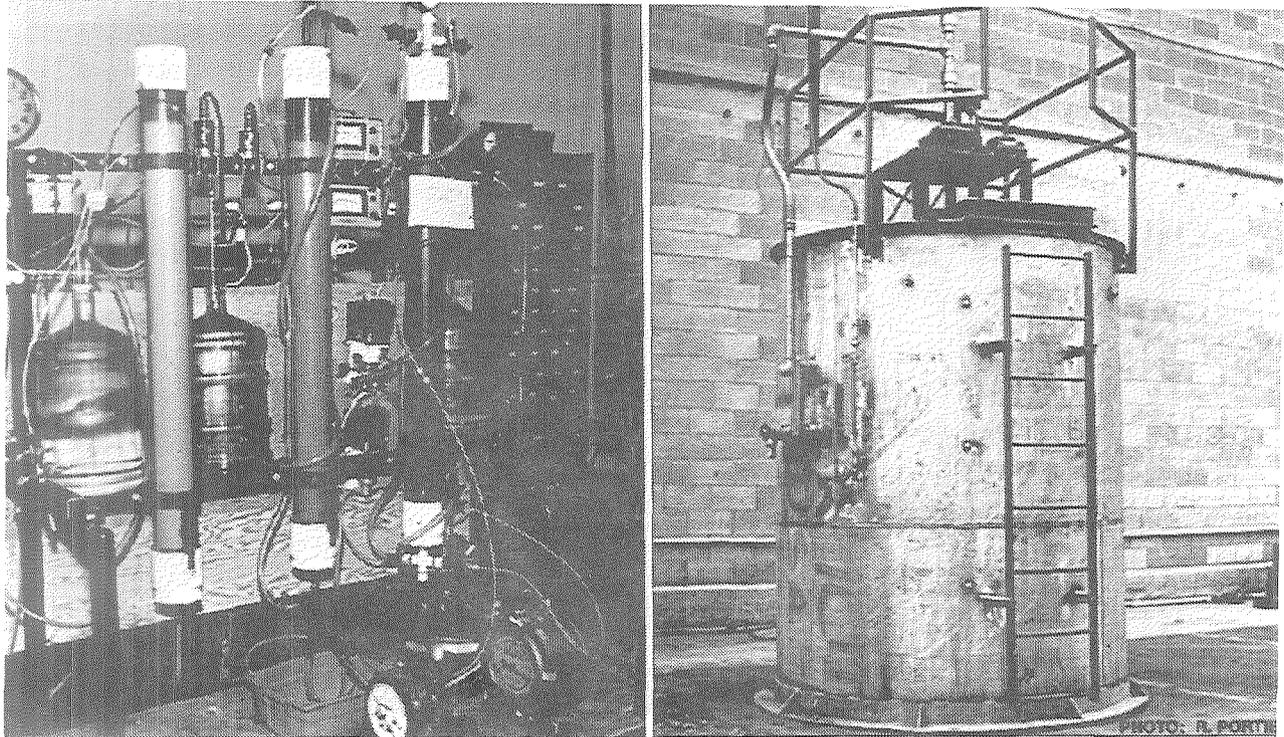


Figure 2. Laboratory bioreactor (left) and field bioreactor for site remediation

One of the methods Lockheed has investigated for processing solid wastes and highly toxic materials is wet oxidation. Wet oxidation is a flameless combustion process carried out at moderate temperatures (300-650°F) and moderate to high pressures (450-2500 PSIG). Pure oxygen gas or compressed air can be used to supply oxygen for the reaction process.

The products of the oxidation process vary, depending on the nature of the waste stream, and the operating conditions of the reactor. Under ideal operating conditions, the products include only CO₂, N₂, H₂O and dissolved inorganic salts. Figure 5 shows a laboratory wet oxidation reactor test bed, along with a mobile reactor test bed and the mobile test bed control room. These reactors were developed and extensively tested in the late 1970's.

Figure 6 illustrates test results obtained by applying the mobile test bed reactor to a variety of plant effluents. As indicated, reduction of total organic carbon in these (TOC) effluents was always at least 44-45%, while the removal of the specific constituent of interest never fell below 88%. These operating efficiencies are controllable by altering the conditions under which the oxidation occurs.

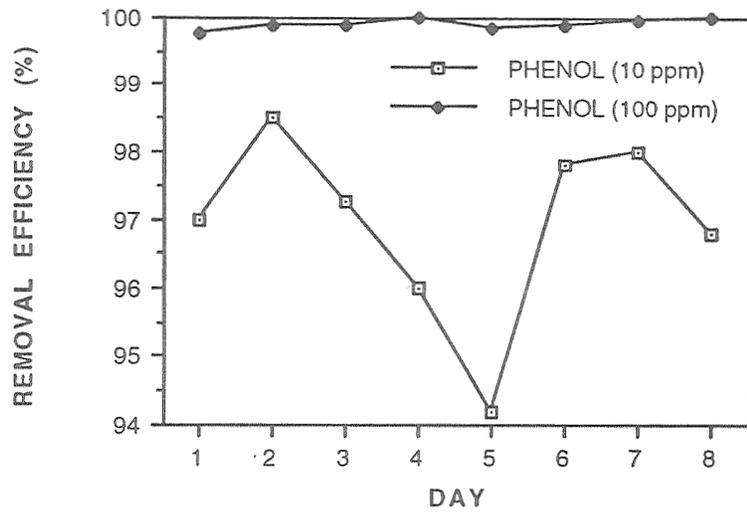


Figure 3. Removal of phenol from water by a microbial bioreactor.

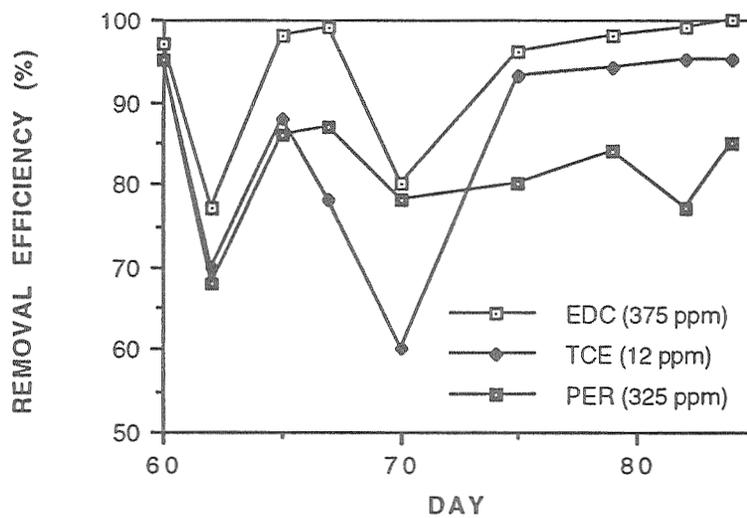


Figure 4. Removal of chlorinated hydrocarbons from water by a microbial bioreactor.

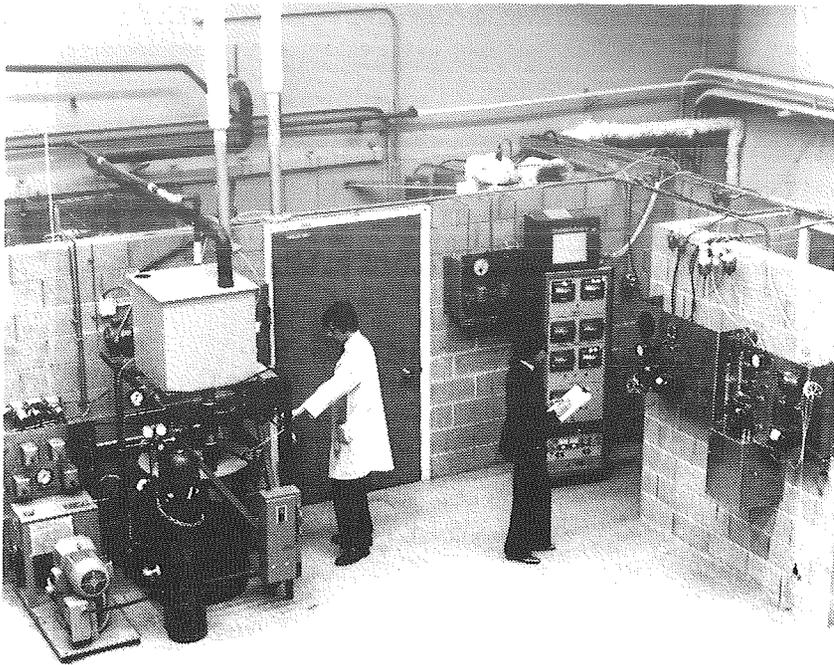


Figure 5. Lockheed wet oxidation reactors: laboratory research test bed, mobile test bed, and mobile test bed control room.

Waste Source (Specific Constituent)	Operating Conditions		TOC Reduction (%)	Specific Constituent Removed (%)
	Temperature (°F)	Pressure (PSIG)		
Benzoic Herbicide Process (Dichlor- nitrobenzoic acid)	500	1,200	45	88
Acrolein Process (Acrolein/allyl alcohol)	550	1,500	44	99
Triazine Herbicide Process (Atrazine derivatives)	500	1,200	46	100
Xylene Process (Aromatics)	600	2,000	68	92
Antiozonant Process (Aromatics)	450	900	64	100
Hydrazine Process (Hydrazine)	400	700	49	100
Urea- Formaldehyde Process (TKN)	250	300	60	92
Coke Plant Ammonia Still	550	1,500	90	100
Synthetic Rubber Process (Surfactants)	600	2,000	67	95

Figure 6. Results obtained from the mobile test bed reactor when applied to a variety of plant effluents.

CONCLUSION

Many of the environmental problems we must solve over the next two decades are life support problems. As a consequence, spacecraft life support technologies, or modifications of those technologies, can provide us with additional methods of solving the problems. In the appropriate configurations, these technologies can contribute substantially to pollution prevention, site restoration, and recycling. Ultimately, these technologies, which are now being developed to support the life of humans on other planets, may play a crucial role in sustaining human life on the planet we call home.

LITERATURE CITED

Miller, G.P., R.J. Portier, and H.L. Sleeper. 1992. Further Applications of the Use of Biological Reactors to remove Trace Hydrocarbon Contaminants from Recycled Water. Paper #921273, 22nd International Conference on Environmental Systems, July 13-16, 1992, Seattle, WA.

Wolverton, B.C., R.C. McDonald, and W.R. Duffer. 1983. Microorganisms and Higher Plants for Waste Water Treatment. *J. Environ. Qual.*, 12(2), page 236.

REPORT DOCUMENTATION PAGE			Form Approved OMB No. 0704-0188
1. AGENCY USE ONLY (Leave blank)		2. REPORT DATE February 1993	3. REPORT TYPE AND DATES COVERED CONFERENCE PUBLICATION
4. TITLE AND SUBTITLE Technology 2002 Volume 1		5. FUNDING NUMBERS	
6. AUTHOR(S) MICHAEL HACKETT, COMPILER			
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) NASA TECHNOLOGY TRANSFER PROGRAM (CODE CU)		8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES) NATIONAL AERONAUTICS AND SPACE ADMINISTRATION WASHINGTON, DC 20546		10. SPONSORING/MONITORING AGENCY REPORT NUMBER NASA CP-3189, VOL. I	
11. SUPPLEMENTARY NOTES			
12a. DISTRIBUTION/AVAILABILITY STATEMENT UNCLASSIFIED - UNLIMITED SUBJECT CATEGORY 99		12b. DISTRIBUTION CODE	
13. ABSTRACT (maximum 200 words) Proceedings from symposia of the Technology 2002 Conference and Exposition, December 1-3, 1992, Baltimore, MD. Volume 2 features 60 papers presented during 30 concurrent sessions.			
14. SUBJECT TERMS technology transfer, computer technology, advanced manufacturing, materials science, biotechnology, electronics		15. NUMBER OF PAGES 552	
		16. PRICE CODE A24	
17. SECURITY CLASSIFICATION OF REPORT UNCLASS	18. SECURITY CLASSIFICATION OF THIS PAGE UNCLASS	19. SECURITY CLASSIFICATION OF ABSTRACT UNCLASS	20. LIMITATION OF ABSTRACT UNLIMITED

*U.S. GOVERNMENT PRINTING OFFICE: 1993-728-150/60020

END DATE JUNE 8, 1993

National Aeronautics and
Space Administration
Code JTT
Washington, D.C.
20546-0001
Official Business
Penalty for Private Use, \$300

SPECIAL FOURTH-CLASS RATE
POSTAGE & FEES PAID
NASA
PERMIT No. G27



POSTMASTER: If Undeliverable (Section 158
Postal Manual) Do Not Return
