

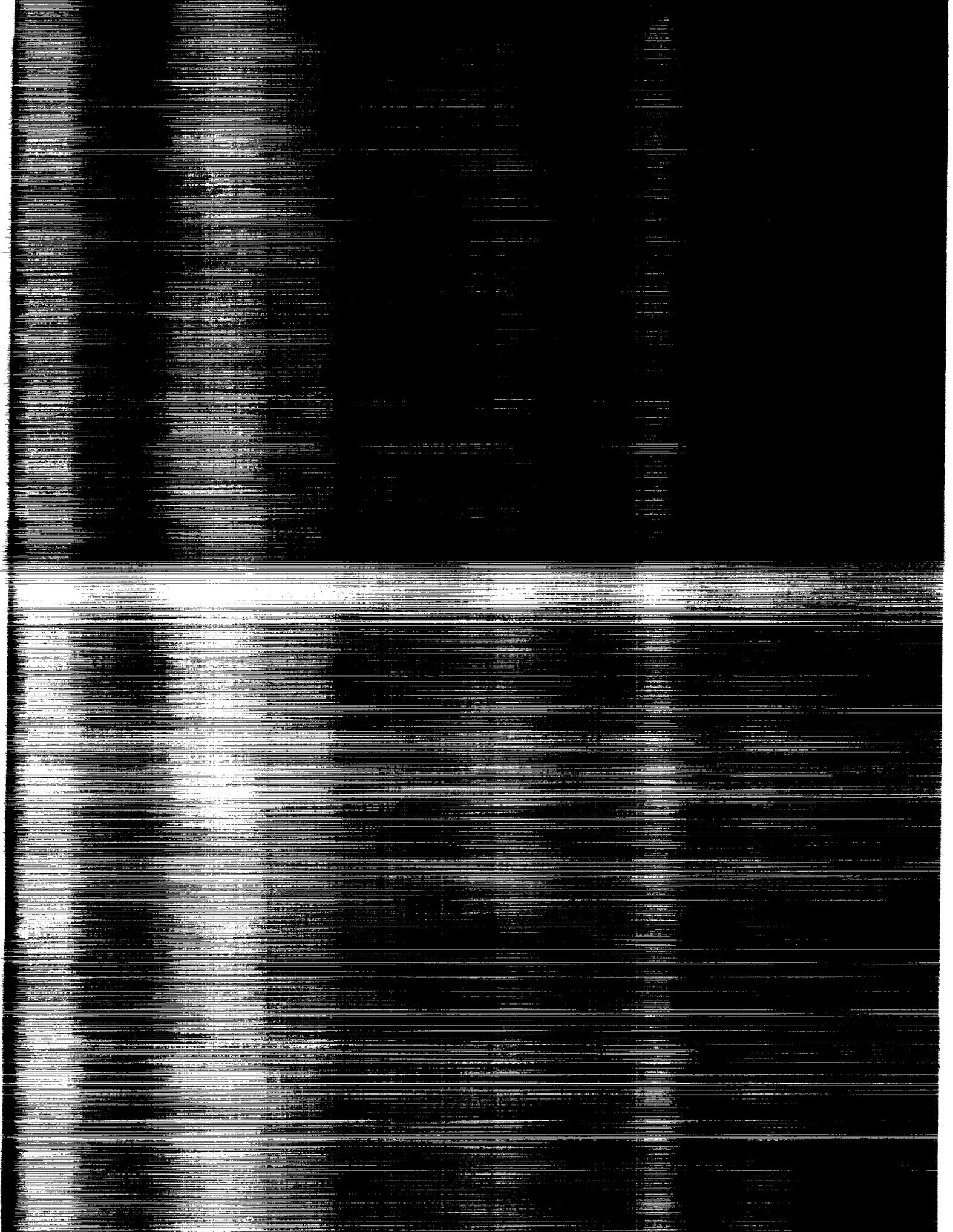
*NASA Conference Publication 3191*

# 1993 Space and Earth Science Data Compression Workshop

(NASA-CP-3191) THE 1993 SPACE AND  
EARTH SCIENCE DATA COMPRESSION  
WORKSHOP (NASA) 119 p

N93-24544  
--THRU--  
N93-24555  
Unclas

H1/59 0152610



*NASA Conference Publication 3191*

# 1993 Space and Earth Science Data Compression Workshop

James C. Tilton, *Editor*  
*NASA Goddard Space Flight Center*  
*Greenbelt, Maryland*

Proceedings of a workshop sponsored by the  
National Aeronautics and Space Administration  
and the IEEE Computer Society Technical Committee  
on Computer Communications and held at  
the Snowbird Conference Center  
Snowbird, Utah  
April 2, 1993

**NASA**

National Aeronautics and  
Space Administration  
Office of Management  
Scientific and Technical  
Information Program

**1993**



## FOREWORD

The third annual Space and Earth Science Data Compression Workshop was held on April 2, 1993 in Snowbird, Utah. This NASA Conference Publication serves as the proceedings for the workshop. The workshop was held in conjunction with the 1993 Data Compression Conference (DCC'93), which was held at the same location March 30 - April 2, 1993.

The goal of the Space and Earth Science Data Compression Workshop series is to explore the opportunities for data compression to enhance the collection and analysis of space and Earth science data. Of particular interest is research that is integrated into, or has the potential to be integrated into, a particular space and/or Earth science data information system. Participants are encouraged to take into account the scientist's data requirements, and the constraints imposed by the data collection, transmission, distribution and archival system.

Papers were selected from submissions to the 1993 Data Compression Conference (DCC '93), and from a limited number of submissions directly to the Workshop. Eleven papers were presented in 4 sessions. Discussion was encouraged by scheduling ample time for each paper, and through scheduled discussion periods at the end of each session.

The workshop was organized by James C. Tilton of the NASA Goddard Space Flight Center, Sam Dolinar of the Jet Propulsion Laboratory, Sherry Chuang of the NASA Ames Research Center, and Dan Glover of the NASA Lewis Research Center. Contact information is given below.

### Acknowledgment

The organization of this workshop was supported by the Office of Advanced Concepts and Technology (and formerly by the Office of Aeronautics and Space Technology), NASA Headquarters, Washington, DC.

### Workshop Organizers

James C. Tilton  
Mail Code 930  
NASA GSFC  
Greenbelt, MD 20771  
phone: (301) 286-9510  
FAX: (301) 286-3221  
Internet:  
tilton@hrpisis.gsfc.nasa.gov  
GSFCMAIL: JTILTON.

Sam Dolinar  
Mail Stop 238-420  
Jet Propulsion Laboratory  
4800 Oak Grove Drive  
Pasadena, CA 91109  
phone: (818) 354-7403  
FAX: (818) 354-6825  
Internet:  
sam@shannon.jpl.nasa.gov

Sherry Chuang  
NASA Ames Research Center  
MS 269-4  
Moffett Field, CA 94035-1000  
phone: (415) 604-3376  
Internet: chuang@ptolemy.arc.nasa.gov

Daniel Glover  
Mail Stop 54-2  
NASA Lewis Research Center  
21000 Brookpark Road, Cleveland, OH 44135  
phone: (216) 433-2847  
Internet: caglove@lims01.lerc.nasa.gov



# CONTENTS

Foreword . . . . .	iii
Contents . . . . .	v

## PAPERS

Data Compression Opportunities in EOSDIS . . . . . <i>Ben Kobler and John Berbert, NASA Goddard Space Flight Center</i>	3
Compression through Decomposition into Browse and Residual Images . . . . . <i>Dmitry A. Novik, Universal Systems and Technology, Inc.; James C. Tilton, NASA Goddard Space Flight Center; and M. Manohar, Universities Space Research Association</i>	7
Applications of Wavelet-Based Compression to Multidimensional Earth Science Data. . . . . <i>Jonathan N. Bradley and Christopher M. Brislawn Los Alamos National Laboratory</i>	13
Wavelet Encoding and Variable Resolution Progressive Transmission . . . . . <i>Ronald P. Blanford, TRW Systems Integration Group</i>	25
Fast Image Decompression for Telebrowsing of Images. . . . . <i>Shaou-Gang Miaou and Julius T. Tou, University of Florida</i>	37
A Survey of Quality Measures for Gray Scale Image Compression . . . . . <i>Ahmet M. Eskicioglu and Paul S. Fisher</i>	49
Digital Mammography, Cancer Screening: Factors Important for Image Compression. . . . . <i>Laurence P. Clarke, G. James Blaine, Kunio Doi, Martin J. Yaffe, Faina Shtern, and G. Stephen Brown, NCI/NASA Working Group on Digital Mammography; Daniel L. Winfield, Research Triangle Institute; and Maria Kallergi, University of South Florida</i>	63
A Study of Video Frame Rate on the Perception of Compressed Dynamic Imagery. . . . . <i>Richard F. Haines and Sherry L. Chuang, NASA Ames Research Center</i>	75
Systems Aspects of COBE Science Data Compression. . . . . <i>I. Freedman, E. Boggess, and E. Seiler, Hughes STX</i>	85
Proposed Data Compression Schemes for the Galileo S-Band Contingency Mission . . . . . <i>Kar-Ming Cheung and Kevin Tong, Jet Propulsion Laboratory</i>	99
Data Compression for the Cassini Radio and Plasma Wave Instrument . . . . . <i>W. M. Farrell, NASA Goddard Space Flight Center; D. A. Gurnett, D. L. Kirchner, and W. S. Kurth, The University of Iowa; and L. J. C. Woolliscroft, The University of Sheffield (United Kingdom)</i>	111



## PAPERS



## Data Compression Opportunities in EOSDIS

Ben Kobler  
 EOS Systems Development Office  
 Mail Code 902.1, NASA GSFC  
 Greenbelt, MD 20771  
 (301) 286-3553  
 (301) 286-3221 (FAX)  
 bkobler@gsofcmail.nasa.gov

John Berbert  
 EOS Systems Development Office  
 Mail Code 902.1, NASA GSFC  
 Greenbelt, MD 20771  
 (301) 286-5916  
 (301) 286-3221 (FAX)  
 jberbert@gsofcmail.nasa.gov

**Abstract.** The Earth Observing System Data and Information System (EOSDIS) is described in terms of its data volume, data rate, and data distribution requirements. Opportunities for data compression in EOSDIS are discussed.

### 1. Introduction

The Earth Observing System Data and Information System (EOSDIS) is being developed by the National Aeronautics and Space Administration (NASA) to be a comprehensive data and information system providing the Earth science research community with easy, affordable, and reliable access to Earth Observing System (EOS) and other appropriate Earth science data. The EOS program, as a part of the Mission to Planet Earth is intended to study global-scale processes that shape and influence the Earth [1, 2, 3]. Beginning in 1998, EOSDIS will archive approximately one terabyte of data per day over a 15 year period [4, 5, 6, 7]. Many opportunities for data compression exist in EOSDIS for alleviating problems due to large data volumes, high bandwidth requirements, and data access requirements.

### 2. EOSDIS Requirements

There are 5 proposed EOS instruments on the EOS AM-1 spacecraft to be launched in June 1998 and 6 proposed EOS instruments on the EOS PM-1 spacecraft to be launched in December 2000. These instruments will generate data at a rate of 281 gigabytes per day [8]. Other instruments will follow on spacecraft to be flown later. Data from the EOS instruments will be transferred to an EOS Data and Operations System (EDOS), from where data will be batched to an appropriate Distributed Active Archive Center (DAAC), selected with responsibility for further processing. The Product Generation System (PGS) located at the DAACs will generate higher level products (L1 through L4) for storage in the Data Archive and Distribution System (DADS). The data product processing levels are defined as follows:

- L0 Raw instrument data at original resolution, time ordered, with duplicate packets removed
- L1A L0 data, which may have been reformatted or transformed reversibly, located to a coordinate system, and packaged with need ancillary and engineering data
- L1B Radio metrically corrected and calibrated data in physical units at full instrument resolution as acquired
- L2 Retrieved environmental variables (e.g. ocean wave height, soil moisture, ice concentration) at the same location and similar resolution as the L-1 source data
- L3 Data or retrieved environmental variables that have been spatially and/or temporally resampled (i.e., derived from L1 or L2 data products) and may include averaging and compositing
- L4 Model output and/or variables derived from lower level data which are not directly measured by the instruments such as new variables based upon time series of L2 or L3 data

Generation of these higher level data products will expand total data volume by a factor of 3.3, resulting in a total data volume from the AM-1 and PM-1 platforms of approximately 0.9 terabytes per day.

The sustained combined daily rate for data input into EOSDIS from the AM-1 and PM-1 platforms will be 26 megabits per second. The sustained daily rate for data access into and from the DADS will, however, be substantially larger to accommodate, in addition to the initial data processing, subsequent data reprocessing and data distribution to users.

A distributed Information Management System (IMS) will be implemented to provide a common user interface to database management systems at the DAACs, providing the capability to easily construct complex queries to search, locate, select, and order products. The IMS will be sized to accommodate 100,000 users. A load of 100 concurrent IMS sessions will be distributed across the DAACs. Approximately 500 IMS queries per hour can be expected for log-on authorization, directory search, catalog search, inventory search, status checks, browse selection, document search, and ordering services.

EOSDIS will be capable of distributing data via physical media and via communications networks, each at a rate equivalent to approximately 1 terabyte per day. Data requested on physical media will be made available for delivery within 24 hours and data requested over networks will be available to the network within an average of 5 minutes.

### **3. Data Compression Opportunities**

Conventional lossless compression techniques such as Huffman coding, Ziv-Lempel compression, and arithmetic coding have been shown to be very effective at compressing a wide range of data types with compression ratios of approximately 2:1 [9, 10, 11]. The potential cost savings to the EOSDIS data archive facility due to reduction of hardware for data storage is obvious. Perhaps less obvious is also a concomitant reduction in requirements for bandwidth of storage devices. To be most effective, however, compressed data needs to stay in its compressed form as long as possible, so that data is not needlessly decompressed and then re-compressed, and so that the potential savings in network bandwidth are not lost. This requires standardization on a common set of data compression schemes, on associated common data format structures, and on common compression/decompression tool kits that are integrated across all of EOSDIS. For example, callable routines that decompress a block or record at a time, would be essential to PGS, as would routines that decompress data at user workstations.

Lossy compression techniques such as DCT, wavelet transform, and vector quantization [12, 13, 14] can play a significant role in optimizing data access by providing tools for storage and retrieval of display quality browse data. EOSDIS will permit users to browse subsetted, subsampled and/or summarized data sets that are created during routine production processing. These browse data sets will be generated by algorithms provided by scientists. Since some of these browse products are designed for visual display, they may be further compressed by lossy compression techniques that can have significantly higher compression ratios than lossless techniques. Because EOSDIS needs to retrieve and begin to display these browse data sets within one minute, they need to be stored on faster access devices than other data. The associated reduction in bandwidth requirements due to data compression could aid in reducing costs.

More innovative lossy/lossless techniques, such as progressive vector quantization [15], have the potential for allowing browse quality lossy compression, while also allowing lossless restoration of full datasets. Such combined techniques can benefit from the best features of both and can result in reduced total I/O requirements and better compression ratios. To be most useful, these techniques require standardization on a common format structure that allows storage of the browse component on a fast access device and storage of the complementary lossless data

component on a slower access device. Unfortunately, however, data compression techniques such as vector quantization are extremely processor intensive, although the decompression phase is much less so. The benefits of reduced I/O and higher compression need to be balanced against the compression cost and the impact of that cost on the PGS design.

Finally, the concept of using very large codebooks to achieve very high compression, both lossless and lossy, although still unproven, has potential for success in extremely large data archives such as those planned in EOSDIS. Fundamental issues need to be investigated that explore the redundancy, and hence compression limit, of these data archives, the stability of the resultant codebooks, and the most effective method for the generation, storage and exchange of those codebooks.

## References

- [1] Ramapriyan, H. K., "The EOS Data and Information System," Proceedings of the AIAA/NASA Second International Symposium on Space Information Systems, Pasadena, CA, September 1990.
- [2] Dozier J. and H. K. Ramapriyan, "Planning for the EOS Data and Information Systems (EOSDIS)," The Science of Global Environmental Change, NASA ASI, 1991.
- [3] Taylor, T. D., H. K. Ramapriyan, and J. Dozier, "The Development of the EOS Data and Information System," Proceedings of the AIAA 29th Aerospace Sciences Meeting, Reno, NV, January 1991.
- [4] Kobler, B. and J. Berbert, "NASA Earth Observing System Data Information System (EOSDIS)," Eleventh IEEE Symposium on Mass Storage Systems, Monterey, CA, 18-19, October 1991.
- [5] Berbert, J. and B. Kobler, "EOSDIS DADS Requirements," NSSDC Conference on Mass Storage Systems and Technologies for Space and Earth Science Applications, Greenbelt, MD, NASA 3165, III-141, July 1991.
- [6] Functional and Performance Requirements Specification for the Earth Observing System Data and Information System (EOSDIS) Core System, NASA, Goddard Space Flight Center, July 1991.
- [7] EOSDIS Core System Statement of Work, NASA, Goddard Space Flight Center, July 1991.
- [8] Lu, Y., "EOS Output Data Products and Input Requirements, Version 2.0," NASA, Goddard Space Flight Center, August 1992.
- [9] Huffman, D. A., "A method for the construction of minimum redundancy codes," Proc. IRE, 40, 1098-1101, 1952.
- [10] Ziv, J. and A. Lempel, "Compression of Individual Sequences via Variable Rate Coding," IEEE Trans. Information Theory, Vol. IT-24, 530-536, September 1978.
- [11] Langdon, G., "An Introduction to Arithmetic Coding," IBM Journal Research Development, Vol. 23, No 2, 135, March 1984.
- [12] Clark, R. J., Transform Coding of Images, Academic Press, London, 1985.
- [13] Antonini, M., M. Barlaud, P. Mathieu, and I. Daubechies, "Image Coding Using Wavelet Transform," IEEE Transactions on Image Processing, Vol. I, No. 2, April 1992.
- [14] Gray, R. M., "Vector quantization," IEEE ASSP Magazine, 1 (2), 4-29, 1984.
- [15] Manohar, M. and J. Tilton, "Progressive Vector Quantization of Multispectral Image Data using a Massively Parallel SIMD Machine," Proceedings of the Data Compression Conference, Snowbird, Utah, 181, March 1992.



## COMPRESSION THROUGH DECOMPOSITION INTO BROWSE AND RESIDUAL IMAGES

Dmitry A. Novik  
Universal Systems and Technology, Inc.  
1000 Wilson Boulevard, Suite 2650  
Arlington, VA 22209  
703-243-7600 x249  
703-243-7788 (FAX)

James C. Tilton  
Information Science and Technology Office  
Mail Code 930  
NASA GSFC  
Greenbelt, MD 20771  
(301) 286-9510  
(301) 286-3221 (FAX)  
tilton@chrpisis.gsfc.nasa.gov

M. Manohar  
Universities Space Research Association  
Mail Code 610.3  
NASA GSFC  
Greenbelt, MD 20771  
(301) 286-3397  
(310) 286-3221 (FAX)  
manohar@chrpalg.gsfc.nasa.gov

**Abstract.** Economical archival and retrieval of image data is becoming increasingly important considering the unprecedented data volumes expected from the Earth Observation System (EOS) instruments. For cost effective browsing the image data (possibly from remote sites), and retrieving the original image data from the data archive, we suggest an integrated image browse and data archive system employing incremental transmission.

We produce our browse image data with the JPEG/DCT lossy compression approach. Image residual data is then obtained by taking the pixel by pixel differences between the original data and the browse image data. We then code the residual data with a form of variable length coding called diagonal coding.

In our experiments, the JPEG/DCT is used at different quality factors ( $Q$ ) to generate the browse and residual data. The algorithm has been tested on band 4 of two Thematic Mapper (TM) data sets. The best overall compression ratios (of about 1.7) were obtained when a quality factor of  $Q=50$  was used to produce browse data at a compression ratios of 10 to 11. At this quality factor the browse image data has virtually no visible distortions for the images tested.

### 1. Introduction

Economical archival and retrieval of image data is becoming increasingly important considering the unprecedented data volumes expected from the Earth Observation System (EOS) instruments. The challenges EOS present to the information scientist are providing a cost effective mechanism for: (i) browsing the image data (possibly from remote sites), and (ii) obtaining the original image data from the data archive. We suggest that these two mechanisms be integrated, *i. e.*, the lossless image data should be reconstructed from the browse image data by incremental transmission.

The data archive's integrity is maintained as long as every bit of the original image data can be reliably reconstructed from compressed form without loss. Nevertheless, lossless compression is not very effective in reducing data volume. Maximum compression ratios of 2.0 to 2.5 are typical for the type of image data expected from EOS instruments. Lossy compression, on the

other hand, can typically provide compression ratios of as high as 30 to 50 without significant visible degradation of the image data. However, because the original image data cannot be perfectly reconstructed from this highly compressed data, it can only be used for data browsing and, possibly, certain preliminary analysis.

In most data archive schemes, highly compressed data is kept in on-line storage and used to efficiently browse the data to determine potentially useful data set(s) for further processing. Once this decision is made, the original data is obtained from off-line storage. The browse quality image data and the corresponding original image data contain redundant information, causing a fraction of the information to be transmitted twice.

If incremental data is stored off-line instead of original data, data transmission to users can be made more efficient. In this approach the image data is decomposed into browse and residue so information is not duplicated either in data archival or in transmission to users across the computer networks.

In this paper we address the problem of decomposing image data into browse and residual data in a manner that is most appropriate for image data archival. Browse data should take only a small fraction (typically 1/30 to 1/50) of the storage required for original data with quality that is adequate for deciding whether the data is useful or not for an intended application. The residual data, normally kept off-line, should have relatively high compressibility using a carefully designed lossless compression technique. Thus, the key problems are to select a lossy compression approach that provides the best compression with quality that is nearly lossless visually, and to select the most effective lossless compression approach for the residual. In addition, we also determine the browse data compression ratio that leads to the best overall compression.

## 2. JPEG/DCT Approach for Browse Quality Image Generation

Any of several lossy compression techniques, such as subband/wavelet coding and vector quantization, could be used to produce the browse quality image. We chose to use the JPEG/DCT lossy compression approach for the following reasons.

- i. The JPEG/DCT lossy compression approach has become an industry-wide standard compression approach.
- ii. Special hardware boards are available commercially for various machines including the ubiquitous IBM/PC.
- iii. The image quality of the browse data can be fine tuned until it is visually lossless.

JPEG lossy compression is based on the Discrete Cosine Transform (DCT) of 8x8 blocks of the input image [1-2]. In the encoding process, the samples in the input image are grouped into 8 x 8 blocks, and each block is transformed by the forward DCT (FDCT) into a set of 64 coefficients referred to as the DCT coefficients. The first coefficient corresponds to the DC coefficient, and the remaining 63 are AC coefficients. Each of the AC coefficients is then quantized using one of 64 corresponding values from a quantization table. The DC coefficients of different blocks undergo differential coding. The AC coefficients are then ordered by a one-dimensional zigzag sequence. Finally, the quantized coefficients are compressed using either a Huffman table or arithmetic coding.

The baseline JPEG/DCT accepts 8-bit images and uses two Huffman tables for coding DC and AC coefficients. However, the other JPEG lossy standards allow 8-bit to 12-bit precisions with either Huffman or arithmetic coding of coefficients. At the decoding end the 64 coefficients are

used to reconstruct 8 x 8 coefficient image which then is mapped back to image space by Inverse DCT (IDCT).

JPEG/DCT approach provides a fine tuning factor,  $Q$ , which corresponds to different qualities of the compressed images. For typical NASA image data, a low value of  $Q$ , such as 20, provides high compression with poor image fidelity. As the  $Q$  factor increases the fidelity improves at the expense of compression ratio. For  $Q = 80$ , the compressed images are generally visually indistinguishable from the input images, with a compression ratio typically in the range of 6.0 to 7.0. For data from the Landsat TM instrument, a general image quality rating for different  $Q$  values and corresponding compression ratios ( $C_R$ ) is:

$Q$	$C_R$	Image Quality
25 - 40	25 - 12	moderate to good quality
40 - 70	12 - 8	good to very good quality
70 - 80	8 - 6	excellent quality
80 - 90	6 - 4	indistinguishable from original

Several EOS instruments are expected to have a dynamic range of 0 - 4095, that is, the pixel brightness level can be represented by 12 bits. However, the human perceptual system cannot even resolve 256 gray scale levels (*i. e.*, a range of 0 - 255), which can be represented by 8 bits. Therefore, the first stage of producing the browse data can be described as follows: Determine the actual dynamic range of the data (which can be less than, but no more than 12 bits), and retain only the 8 most significant bits in that dynamic range. Then compress this 8-bit data with JPEG/DCT at the optimal quality factor. For lossless compression, the remaining bits, as well as the residual from JPEG compression are separately compressed using an appropriate lossless compression approach. Such an approach is described in the following section.

### 3. Residual Compression using Diagonal Codes

Residual image data is that which is obtained through taking the pixel by pixel differences between the original data and the image reconstructed after lossy compression. We have observed that the residual image data obtained from JPEG/DCT compression is low entropy data that is compressible to a greater degree than the original image data. The better the browse data approximates the original data, the more compressible is the residual image data. Thus, a better quality browse results in a residual that can be compressed better in lossless mode.

However, a better quality browse image requires more bits per pixel. Since the overall lossless representation is sum of the bits per pixel for browse data and residual data, producing maximum overall lossless compression requires finding the optimal balance between the bits allocated to the browse data and the bits consequently required for the residual data.

For remote browsing applications, the browse data bit rate (bits/pixel) must be kept very low to ensure efficient transmission of the data across the computer networks. This requirement leads to choosing the lowest JPEG/DCT quality factor without significant visual degradation of the reconstructed image data, which we have found to be a quality factor of about 50. Fortunately, our experiments have found that a quality factor of about 50 also corresponds closely to the browse bit rate that produces the optimal overall lossless compression in combination with the residual image data.

The residual image data exhibits a Laplacian distribution with a smaller variance of data values than the original image data. This property suggests that a form of variable length encoding

would be most appropriate for lossless compression of this data. We have found specifically that a type of variable length encoding, called diagonal coding [3,4], is most appropriate.

For images with  $n$  bits/pixel, straightforward representation of the residual image data requires  $n+1$  bits. However, through using Golomb codes [5], the residual data requires just  $n$  bits/pixels (prior to diagonal coding).

In our approach, the residual image data is divided into two parts. The first part contains the lower order two bits, while the second part contains the remaining higher order six bits. The frequency distribution of the lower order bits exhibits no particular structure, and thus can be compressed very little. However, the frequency distribution of the higher order bits exhibits a narrow Laplacian distribution. For this type of distribution, Rice, *et. al.*, [3] have shown the diagonal code is asymptotically optimal. In this code, each value is represented by number of zeros corresponding to that value, terminated by a one. For six bit data, the diagonal code for "000101" is "000001", and the diagonal code for "010100" is "00000000000000000001." Since higher values in the residual data occur less frequently, this code turns out to be optimal. This representation is very efficient for coding as well as decoding.

The diagonal code we propose is as follows. The frequency distribution is divided into sets of four pixels centered about the zero axis such that each set contains two negative and two positive residual values except the first one that contains zero. If the residual value belongs to set 1, it is represented by 1, if the value belonged to set 2, it is represented by 01, if the value belongs to set 3, it is represented by 001, and so on. In general, if the residual value belongs to  $i^{\text{th}}$  set, the representation is series of  $i-1$  zeros followed by 1. Typical sets and their representations are shown below:

<u>Set</u>	<u>Range</u>	<u>Diagonal code</u>
1	(-1,0,1,2)	= 1 followed by two bits for identification of actual value
2	(-2,-3,2,4)	= 01
3	(-5,-4,5,6)	= 001
4	(-7,-6,7,8)	= 0001
5	(-9,-8,9,10)	= 00001
6	(-11,-10,11,12)	= 000001
7	(-13,-12,13,14)	= 0000001
8	(-15,-14,15,16)	= 00000001
9	(-17,-16,17,18)	= 000000001
10	(-19,-18,19,20)	= 0000000001
11	(-21,-20,21,22)	= 00000000001
12	(-23,-22,23,24)	= 000000000001
13	(-25,-24,25,26)	= 0000000000001
14	(-27,-26,27,28)	= 00000000000001
15	(-29,-28,29,30)	= 000000000000001
16	(-31,-30,31,32)	= 0000000000000001

#### 4. Experimental Results and Conclusions

We have tested our compression approach on band 4 of Landsat Thematic Mapper Images of Washington, DC and of Davidsonville, LA (northwest of New Orleans, LA). The browse data was generated using JPEG/DCT at quality factors of 25, 50, and 75. Table 1 shows the frequency distribution of residual data at these quality factors:

Table 1. Washington, DC residual data image statistics.

Diagonal Code Set #	<u>Q = 25</u>	<u>Q = 50</u>	<u>Q = 75</u>
1	.3393	.4101	.4911
2	.2556	.2883	.3055
3	.1794	.1686	.1381
4	.1102	.0819	.0487
5	.0599	.0339	.0127
6	.0311	.0118	.0030
7	.0138	.0036	.0005
8	.0060	.0010	.0001
9	.0026	.0002	.00004

The compression performance of the algorithm is summarized in Tables 2 and 3 for the two data sets we have used in our experiments. For three different Quality factors, the browse compression ration ( $CR_B$ ), the overall lossless compression (CR), and the ratio of CR to the first order entropy ( $CR_e$ ) are tabulated. From the table we see that the best compression ratio in lossless mode corresponds to a quality factor,  $Q = 50$ .

Table 2. Washington D.C. (Band 4)

<u>Q</u>	<u><math>CR_B</math></u>	<u>CR</u>	<u><math>CR/CR_e</math></u>
25	20.0	1.63	0.972
50	11.6	1.67	0.985
70	7.5	1.64	0.977

Table 3. New Orleans (Band 4)

<u>Q</u>	<u><math>CR_B</math></u>	<u>CR</u>	<u><math>CR/CR_e</math></u>
25	16.1	1.599	.9198
50	10.3	1.653	.9901
75	6.9	1.633	.9774

We have described a method of decomposing image data into a browse image and residual image data for active archival and distribution of data. We have found that a variant of diagonal code proposed by us gives the best compression ratio for a residual corresponding to the browse data generated by JPEG/DCT for a quality factor of 50. This quality factor provides browse quality that has very little visible distortions for the images tested.

### References

- [1] Pennebaker, "JPEG Technical Specification, Revision 8," *Working Document No. JTC1/SC2/WG10/JPEG-8-R8* (Aug. 1990).
- [2] Wallace, G. K., "The JPEG still Picture Compression Standard," *Communications of the ACM*, Vol. 34, No. 4, April 1991, pp.31-91.

- [3] Rice, R. F., Yeh, P.-S. and Miller, W. H., " Algorithms for a very high speed universal noiseless coding module," *JPL Publication 91-1*, Jet Propulsion Laboratory, Pasadena, California, Feb. 15, 1991.
- [4] Novik, D. A., *Efficient Coding* (in Russian), published by Energia, Moscow-Leningrad, 1965, p. 46.
- [5] Golomb, S. W., "Run-Length Encodings," *IEEE Trans. on Information Theory*, Vol. IT-12, 1966, pp. 399-401.

# Applications of Wavelet-Based Compression to Multidimensional Earth Science Data

Jonathan N. Bradley and Christopher M. Brislawn  
Computer Research Group, Los Alamos National Laboratory  
Mail Stop B-265, Los Alamos, New Mexico 87545

## ABSTRACT

A data compression algorithm involving vector quantization (VQ) and the discrete wavelet transform (DWT) is applied to two different types of multidimensional digital earth-science data. The algorithm (WVQ) is optimized for each particular application through an optimization procedure that assigns VQ parameters to the wavelet transform subbands subject to constraints on compression ratio and encoding complexity.

Preliminary results of compressing global ocean model data generated on a Thinking Machines CM-200 supercomputer are presented. The WVQ scheme is used in both a predictive and nonpredictive mode. Parameters generated by the optimization algorithm are reported, as are signal-to-noise ratio (SNR) measurements of actual quantized data. The problem of extrapolating hydrodynamic variables across the continental landmasses in order to compute the DWT on a rectangular grid is discussed.

Results are also presented for compressing Landsat TM 7-band data using the WVQ scheme. The formulation of the optimization problem is presented along with SNR measurements of actual quantized data. Postprocessing applications are considered in which the seven spectral bands are clustered into 256 clusters using a  $k$ -means algorithm and analyzed using the Los Alamos multispectral data analysis program, SPECTRUM, both before and after being compressed using the WVQ program.

## I. INTRODUCTION.

This work describes the application of an image compression algorithm involving the discrete wavelet transform and vector quantization to two problems involving earth science data. The coding of outputs of supercomputer-generated global climate model (GCM) ocean simulations and Landsat Thematic Mapper (TM) multispectral imagery is investigated. The compression algorithm has its origins in the coding of gray-scale imagery [1, 2]. A set of vector quantizers is designed (one for each subband in the wavelet decomposition) with parameters selected from the solution of an optimization problem that is formulated to minimize quantization distortion with constraints on the overall bit rate and encoding complexity. Although both data types considered in this work are of dimensionality higher than two, we restrict the discussion to coding implementations based on 2-D transforms.

Compression of the GCM data is approached by both a straightforward two-dimensional extension of the earlier algorithm and a predictive scheme in which two-dimensional prediction residuals are coded. The Landsat images are coded by modifying the bit allocation algorithm to allocate coder resources simultaneously among all of the spectral components. For all scenarios, measurements of quantization distortion are presented as functions of compression ratio and encoding complexity, thus revealing tradeoffs involved in the system design.

## II. MULTIDIMENSIONAL WAVELET TRANSFORM-VECTOR QUANTIZATION.

The data-coding technique used in this work, known as the *wavelet-vector quantization* (WVQ) algorithm, is based on vector quantization of the subbands resulting from a discrete wavelet transform (DWT) decomposition of the data signal. For signals in two or more dimensions, the transform used is based on product filter banks (i.e., tensor products of one-dimensional DWT filters). A  $d$ -dimensional signal transformed in this manner with a two-channel filter bank yields  $2^d$  subbands, any of which can be cascaded back through the filter bank to produce a multirate decomposition of the original signal. Although we are currently working on three-dimensional wavelet transforms for use with the three-dimensional climate model data under investigation, the DWT results presented here are restricted to the case of two-dimensional data fields.

Single-level 2-D DWT analysis and synthesis filter banks are depicted in Figures 1 and 2. The analysis filters ( $H_i$ ) and synthesis filters ( $F_i$ ) used in this paper are biorthogonal linear phase FIR wavelet filters constructed in [3, 4]. Note the use of binary subscripts on the subbands,  $a_{ij}$ , to indicate the filters applied to the rows and columns of the signal,  $x$ . Signals obtained from sampling smooth, continuous data fields usually have most of their energy (or variance) concentrated in the low-frequency part of the spectrum, so it is usually most efficient to cascade only the lowpass-lowpass filtered subband,  $a_{00}$ , back through the analysis bank in Figure 1; the resulting subband is denoted  $a_{00,00}$ , or  $a_{,00}$  for short. This cascade is typically carried down four levels (or so), to the  $a_{,,,00}$  band, which will then contain a large portion of the signal energy concentrated in a heavily downsampled signal component. Note that the downsample factor for an  $l^{\text{th}}$ -level subband in a  $d$ -dimensional scheme using an  $M$ -channel coder is  $m_i = M^{ld}$ , so, e.g., subband  $a_{,,,00}$  in the 2-D decomposition has been downsampled by  $m_i = 256$ -to-1.

A further consideration when applying a DWT to finite-duration signals, like the rows or columns in a digital image, is the handling of boundary conditions. The most straightforward way of dealing with signal boundaries is to regard the signal as a single period of an infinite, periodic input and apply the DWT filter bank by circular convolution and downsampling. This has the disadvantages, however, of introducing a spurious jump discontinuity when the data isn't inherently periodic and of constraining the signal length (i.e., its "period") to be divisible by the downsample factor. For a four-level decomposition using a two-channel filter bank, for instance, this means the input length,  $N_0$ , must be divisible by 16. Both of these problems can be avoided by using symmetric extrapolation techniques to extend finite-duration inputs; moreover, this can be done with no increase in the memory allocation needed to transform or store the input signal [5, 6].

The design of vector quantizers for the subbands in a DWT decomposition is based on

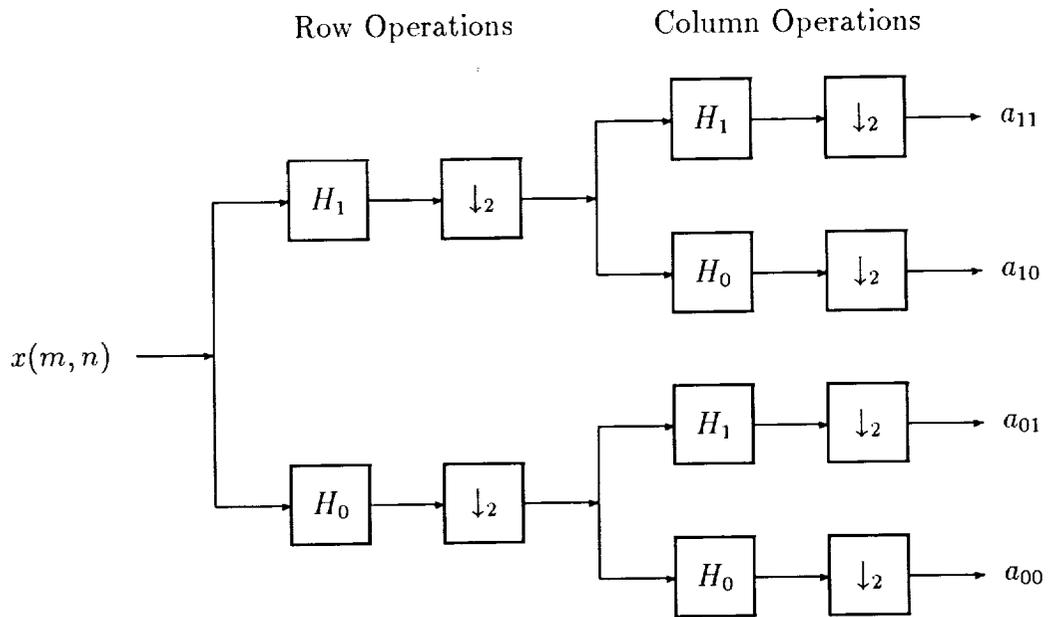


Figure 1: Single-Level, Two-Dimensional Wavelet Transform Analysis.

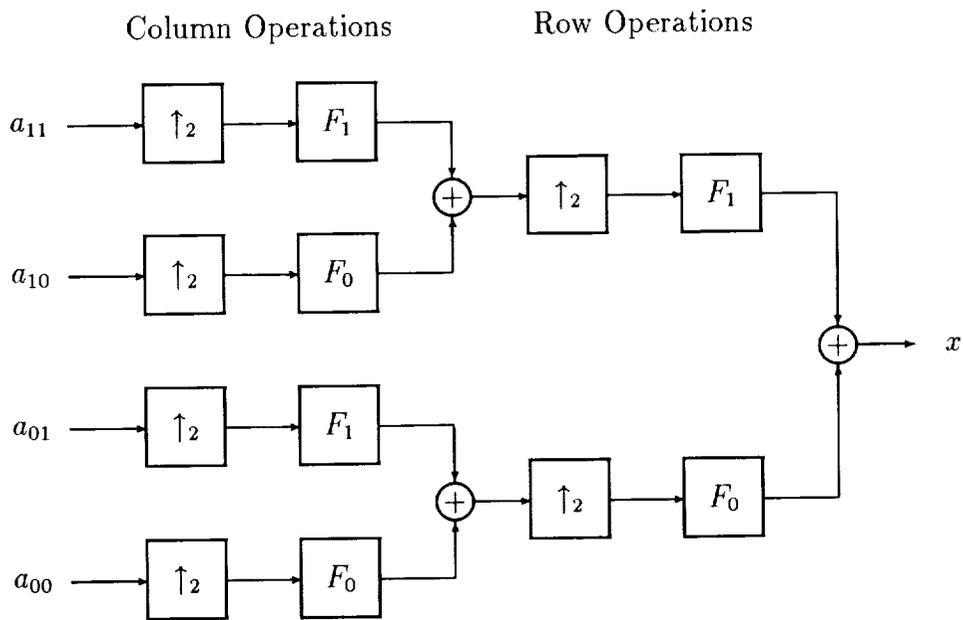


Figure 2: Single-level, Two-Dimensional Wavelet Transform Synthesis.

the assumption of exponential VQ rate-distortion characteristics,

$$D_i(k_i, r_i) = \beta_i(k_i)e^{-\gamma_i(k_i)r_i} \quad (1)$$

$D_i(k_i, r_i)$  is the distortion (mean-square error) between the original and quantized data in the  $i^{\text{th}}$  subband for a bit rate of  $r_i$  bits per pixel (bpp) and a vector dimension of  $k_i$ .  $\beta_i(k_i)$  and  $\gamma_i(k_i)$  are constants that depend on  $k_i$  and the probability density function of the data vectors. The motivation for this assumption is based on theoretical VQ rate-distortion modelling [7] and confirmed by empirical data; values for the constants  $\beta_i(k_i)$  and  $\gamma_i(k_i)$  are determined from a set of training data.

In the case of an orthogonal subband decomposition, the overall distortion can be expressed as a weighted sum of the distortion in each subband,

$$D = \sum_i \frac{1}{m_i} D_i(k_i, r_i) \quad ,$$

where  $m_i$  is the downsample factor, the ratio of the number of samples in the original to the number in the  $i^{\text{th}}$  subband. Since the DWT conserves the number of data samples,  $m_i$  satisfies the identity  $\sum m_i^{-1} = 1$ . By (1), the overall distortion is

$$D = \sum_i \frac{1}{m_i} \beta_i(k_i) e^{-\gamma_i(k_i)r_i} \quad . \quad (2)$$

Formula (2) is customarily used as a distortion measure with nonorthogonal transforms, too, although it no longer coincides exactly with overall mean-square error. The bit-allocation problem for quantizer design involves using nonlinear optimization techniques to compute the bit rates,  $r_i$ , and dimensions,  $k_i$ , that minimize (2), subject to constraints on overall bit rate and encoder complexity.

For a target overall bit rate of  $R$  bpp, the constraint on subband bit rates is

$$\sum_i \frac{r_i}{m_i} \leq R \quad . \quad (3)$$

If subband vector dimensions,  $k_i$ , are to be optimized, an additional constraint besides (3) is necessary to obtain a well-posed optimization problem for VQ bit allocation. The encoder complexity constraint used here is an upper bound,  $Q$ , on the computational cost of performing exhaustive nearest-neighbor searches of  $k_i$ -dimensional VQ codebooks containing  $N_i = 2^{k_i r_i}$  codevectors:

$$\sum_i \frac{1}{m_i} 2^{k_i r_i} \alpha \leq Q \quad . \quad (4)$$

The parameter  $\alpha$  is a constant corresponding to the arithmetic cost of performing two additions and one multiplication per pixel. With the additional constraints  $r_i \geq 0$  we obtain a convex nonlinear optimization problem to solve for the  $r_i$ ; the  $k_i$  are optimized by a heuristic search procedure. Once optimal bit rates and vector dimensions are computed, optimal VQ codebooks are constructed from training data using the Linde-Buzo-Gray method [8, 9]. More details about the WVQ algorithm are given in [10, 2, 1, 11].

### III. APPLICATION TO OCEAN MODEL DATA.

This section describes the use of the WVQ algorithm on synthetic data generated by a Bryan-Cox-Semtner global ocean circulation model running on the Connection Machine CM-200 at the Los Alamos Advanced Computing Laboratory (ACL) [12, 13]. The model is computed on a  $320 \times 768$  grid at 20 depth levels; boundary conditions are given on a three-dimensional bottom topography with 80 islands. The data used in the compression experiments was the surface temperature field (no depth components), taken at three-day intervals over a decade's worth of simulation. Time-frames from the first year of the simulation were used for training data, and the resulting WVQ algorithm was then tested on frames from the last year of the simulation, i.e., on data similar but not identical to the training data. We feel this is a valid test since it is similar to the manner in which the algorithm will be used in practice.

The two-dimensional data frames were transformed with a four-level octave-scaled DWT decomposition. Since the model is periodic in the east-west direction, periodic boundary conditions were used along parallels of latitude ("rows"). However, due to the lack of continuity between the north and south edges of the grid, symmetric (i.e., reflected) boundary conditions were used along meridians of longitude ("columns"). Because the DWT is most easily computed on a rectangular grid, the temperature data was extended across the continental landmasses before transforming. A simple approach like zero-padding of the data would be undesirable because it would induce a large jump discontinuity around the coastlines; this would show up as added variance in the highpass-filtered DWT subbands and would therefore reduce the compressibility of the high-frequency signal components in the transform domain. For this reason we used a continuous extension of the data given by linear extrapolation from coast to coast along parallels of latitude. This still leaves a "corner" at the coastlines in the extended data; since the initial data field is extremely smooth, this corner results in a slight increase in energy in highpass-filtered subbands. It is not yet clear whether this added variance is significant alongside the variance naturally present in the data. We are currently looking into using smoother two-dimensional extrapolation schemes for this task.

Two different approaches were taken to quantizing the time-series data generated in the simulation: nonpredictive and predictive coding. In the nonpredictive scheme, each frame is treated as a separate image and compressed accordingly using the WVQ method. In predictive coding, a prediction of each frame is made based on past frames and subtracted from the current frame, resulting in a two-dimensional *residual* image, which is compressed and stored. The image sequence is decoded from the first frame in the sequence and the residuals. We used a simple first-order predictor in this scheme; i.e., the prediction of a given frame is just equal to the quantized value of the previous frame. Block diagrams for the transmitter and receiver in this predictive encoding/decoding system are given in Figures 3 and 4. The experiment assumed that the first frame in the sequence is transmitted with nonpredictive quantization, and the compression ratios reported are those of the subsequent residuals.

For both the nonpredictive and predictive schemes, WVQ coders were designed for bit rates,  $R$ , ranging from 2.0 to 0.25 bpp. Since the original data was 32 bpp, the corresponding compression ratios range from 16:1 to 128:1. Encoding complexities,  $Q$ , were varied between

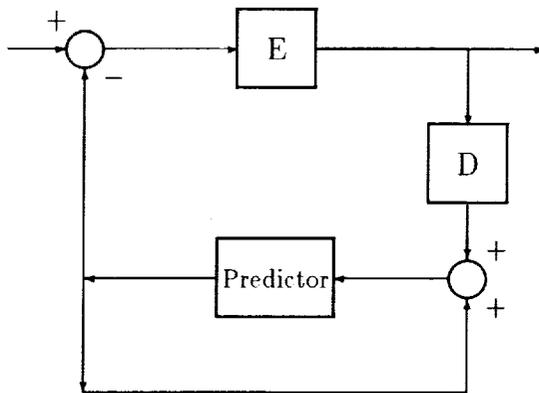


Figure 3: Predictive Transmitter.

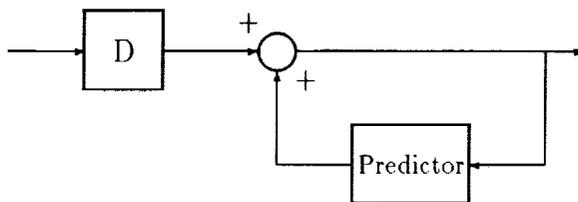


Figure 4: Predictive Receiver.

$16\alpha$  and  $64\alpha$ . The optimal codebook sizes and vector dimensions for each combination of  $R$  and  $Q$  were computed by the WVQ design algorithm described in Section II. Bit allocation results for the 13 subbands are presented in Tables I and II for  $R = 0.5$  and  $0.25$  bpp in terms of vector dimensions,  $k_i$ , and codebook sizes,  $N_i$ . Note that as the bit rate decreases, the highest frequency subbands are quantized more heavily or discarded altogether and remaining high-frequency subband vector dimensions typically increase. Vector dimensions also increase as the upper bound on complexity increases. Bit allocation for the residual subbands in Table II is similar to that for the nonpredictive scheme, although much less of the quantizer resources (i.e., far fewer bits) are allocated to the lower frequency subbands in the predictive scheme. This means that the first-order predictor effectively predicts the low-wavenumber modes of the model, indicating that these modes are evolving slowly compared to the sampling rate for archiving data.

Quantizer performance is measured in terms of signal-to-noise ratio (SNR),

$$\text{SNR} = 10 \log_{10} \frac{\sigma_s^2}{\sigma_e^2} \quad (\text{dB}) \quad ,$$

where  $\sigma_s^2$  is the signal variance and  $\sigma_e^2$  is the quantization error variance. The average SNR in dB for the test data is shown in Figures 5 and 6 for various values of  $R$  and  $Q$ . This diagram illustrates the various tradeoffs involved in the selection of  $R$  and  $Q$ . At a given bit rate,  $R$ , note that a higher SNR is possible using an encoder with higher complexity,  $Q$ ; i.e., higher subband vector dimensions,  $k_i$ , and correspondingly larger codebook sizes,  $N_i$ . Note

Table I: Vector Dimension and Codebook Size Assignments ( $k, N$ ) for  $R = 0.5$  bpp and  $R = 0.25$  bpp, Nonpredictive Coding.

	$R = 0.5$			$R = 0.25$		
	$Q = 64\alpha$	$Q = 32\alpha$	$Q = 16\alpha$	$Q = 64\alpha$	$Q = 32\alpha$	$Q = 16\alpha$
Subband ,,,00	(1,511)	(1,493)	(1,308)	(1,251)	(1,251)	(1,251)
Subband ,,,01	(1,132)	(1,119)	(1,91)	(1,32)	(1,31)	(1,32)
Subband ,,,10	(1,93)	(1,103)	(1,75)	(1,15)	(1,17)	(1,18)
Subband ,,,11	(1,48)	(1,58)	(1,41)	(1,8)	(1,7)	(1,7)
Subband ,,01	(4,723)	(2,247)	(2,133)	(4,377)	(4,382)	(4,295)
Subband ,,10	(4,261)	(2,61)	(2,41)	(4,16)	(4,16)	(4,17)
Subband ,,11	(4,134)	(4,76)	(2,25)	(4,4)	(4,4)	(4,4)
Subband ,01	(8,373)	(8,177)	(4,68)	(8,176)	(8,176)	(8,132)
Subband ,10	(8,19)	(8,13)	(8,8)	—	—	—
Subband ,11	—	—	(8,2)	—	—	—
Subband 01	(8,74)	(8,42)	(8,22)	(8,4)	(8,4)	(8,4)
Subband 10	—	—	—	—	—	—
Subband 11	—	—	—	—	—	—

Table II: Vector Dimension and Codebook Size Assignments ( $k, N$ ) for  $R = 0.5$  bpp and  $R = 0.25$  bpp, Predictive Coding.

	$R = 0.5$			$R = 0.25$		
	$Q = 64\alpha$	$Q = 32\alpha$	$Q = 16\alpha$	$Q = 64\alpha$	$Q = 32\alpha$	$Q = 16\alpha$
Subband ,,,00	(1,14)	(1,10)	(1,16)	(1,2)	(1,2)	(1,4)
Subband ,,,01	(1,21)	(1,15)	(1,23)	(1,3)	(1,4)	(1,6)
Subband ,,,10	(1,26)	(1,20)	(1,28)	(1,5)	(1,6)	(1,9)
Subband ,,,11	(1,35)	(1,28)	(1,35)	(1,8)	(1,10)	(1,14)
Subband ,,01	(4,383)	(4,235)	(2,84)	(4,207)	(4,198)	(4,120)
Subband ,,10	(4,423)	(2,125)	(2,86)	(4,293)	(4,257)	(4,140)
Subband ,,11	(4,305)	(4,185)	(2,64)	(4,93)	(4,109)	(4,87)
Subband ,01	(8,276)	(4,126)	(4,62)	(8,461)	(8,209)	(8,78)
Subband ,10	(8,239)	(8,138)	(8,54)	(8,206)	(8,145)	(8,65)
Subband ,11	(8,97)	(8,51)	(8,24)	—	(8,3)	(8,12)
Subband 01	(8,30)	(8,12)	(8,11)	—	—	—
Subband 10	—	—	—	—	—	—
Subband 11	—	—	—	—	—	—

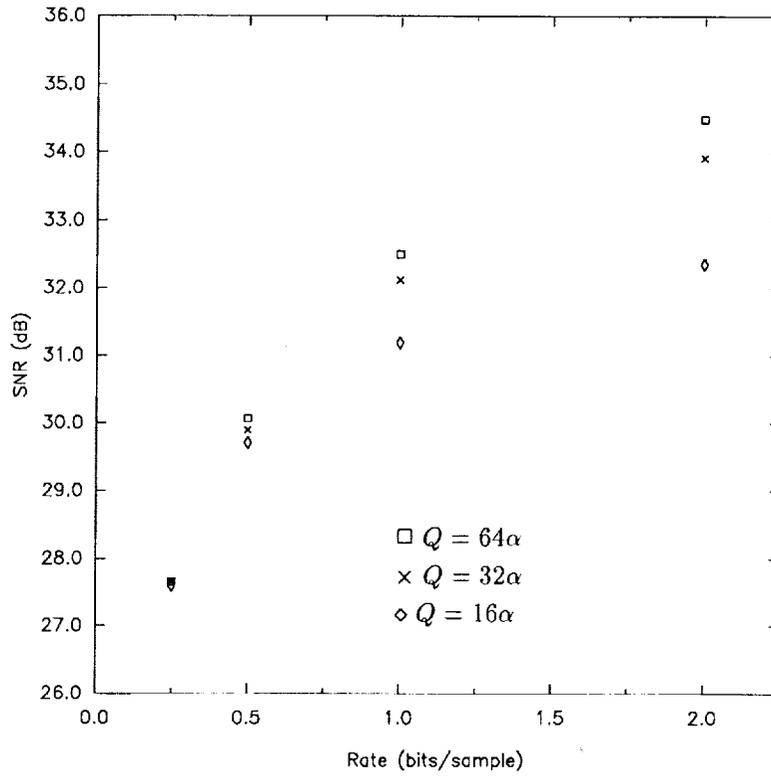


Figure 5: SNR Measurements of Quantized Temperature Data, Nonpredictive.

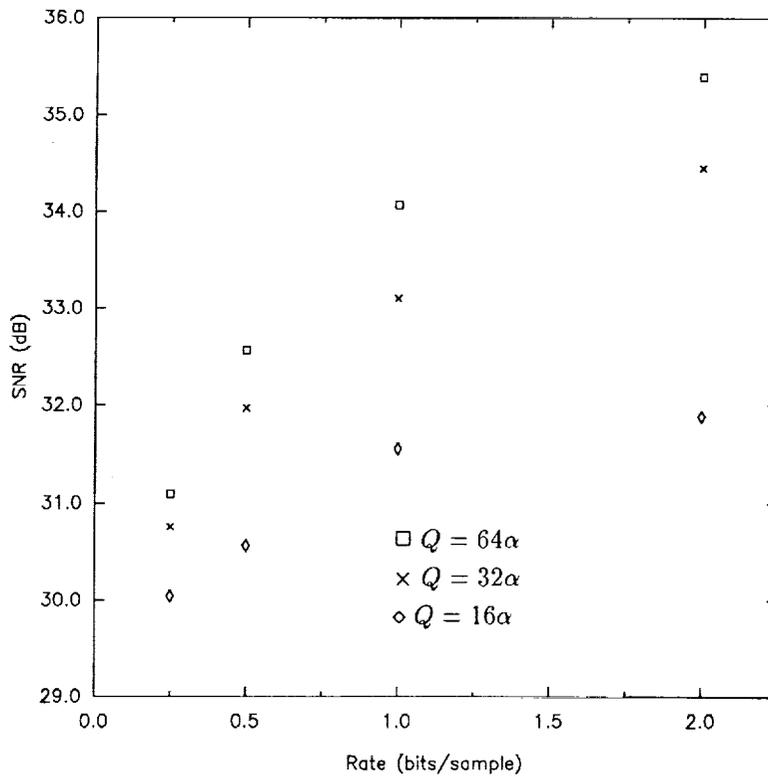


Figure 6: SNR Measurements of Quantized Temperature Data, Predictive Coding.

that the complexity can be increased in this manner while the subband bit rates,  $r_i$ , remain unchanged since

$$r_i = \frac{\log_2 N_i}{k_i} .$$

Increasing the encoding complexity results in an encoder with a more time-consuming code-book lookup but does not affect decoder performance. As the bit rate increases, we see from Figure 5 that the gain in SNR achieved by increasing the encoding complexity becomes more significant. For a fixed  $R$  and  $Q$ , comparison of Figures 5 and 6 shows that the predictive scheme achieves a gain on the order of 1–4 dB over nonpredictive coding. The improvement in coding gain is more pronounced at lower bit rates and higher complexities, since the residual coding scheme is better able to exploit higher limits on encoding complexity at low bit rates than the nonpredictive scheme.

#### IV. APPLICATION TO MULTISPECTRAL DATA.

This section discusses the application of WVQ to the compression of multispectral imagery. Since each spectral band is a separate monochromatic image, the approach is to code each of the bands by two-dimensional WVQ using symmetric boundary conditions. The bit-allocation is performed for the various spectral components simultaneously and hence the coding of each spectral component is not viewed as a separate two-dimensional problem.

The multispectral problem requires a modification to the WVQ design algorithm discussed in Section II since the rate and complexity are expressed in terms of multidimensional pixels. For the case of  $L$  spectral components the system design procedure entails minimizing

$$D = \frac{1}{L} \sum_i \frac{1}{m_i} \beta_i(k_i) e^{-\gamma_i(k_i)r_i} \quad (5)$$

over the  $k_i$  and  $r_i$  subject to

$$\frac{1}{L} \sum_i \frac{r_i}{m_i} \leq R \quad (6)$$

$$\frac{1}{L} \sum_i \frac{1}{m_i} 2^{r_i k_i} \alpha \leq Q \quad (7)$$

$$r_i \geq 0 \quad (8)$$

$$k_i \in K_i \quad (9)$$

where  $K_i$  denotes a prespecified set from which  $k_i$  must be selected. The optimization is performed over all of the two-dimensional subbands generated from all of the spectral components.

Multispectral image WVQ was considered for the application of compressing Landsat Thematic Mapper (TM) data. Such data consist of seven 8-bit spectral bands (three visible, three infrared, and one thermal) at a ground sample distance of 28.5 meters. Four data sets were used in training the coder: Albuquerque, NM ( $2984 \times 3356$ ); Cairo, Egypt ( $2945 \times 3320$ ); Los Alamos, NM ( $2984 \times 3254$ ); and Mexico City, Mexico ( $5965 \times 6967$ ). The performance of the coder was evaluated in terms of results obtained by compressing a ( $2976 \times 3552$ ) scene from the Moscow, Russia, area containing both urban and agricultural areas. The resulting

Table III: RMSE Quantizer Performance as a Function of Compression Ratio and Complexity.

	$Q = 8\alpha$	$Q = 16\alpha$	$Q = 32\alpha$	$Q = 64\alpha$
16:1	3.04	2.35	1.92	1.70
32:1	3.06	2.63	2.43	2.26
64:1	3.85	3.51	3.32	3.24
128:1	4.76	4.71	4.58	4.56

root mean-square error (RMSE) for sixteen combinations of bit rate and encoding complexity are tabulated in Table III. The compression ratios reported are relative to 56 bits/pixel in the original data and assume that the bit rates satisfy

$$\frac{1}{L} \sum_i \frac{r_i}{m_i} = R \quad , \quad \text{with } L = 7 \quad . \quad (10)$$

The additional gain available from entropy and run-length coding is not included.

It is interesting to compare these results to those obtained by another wavelet-based compression technique. In [14] Landsat TM images were compressed via a subband decomposition of each spectral component by a 7-tap nonperfect reconstruction filter bank; each subband was coded with uniform scalar quantization followed by Huffman and zero-run-length coding. The experiment was repeated with an image-dependent Karhunen-Loeve transform (KLT) in the interband direction, which provided noticeable coding gain at the expense of computational complexity. The WVQ RMSE results depicted in Table III appear to lie between these two previous approaches, although any comparisons must be qualified by the fact that the numerical results in these two papers were obtained from different imagery (Kuwaiti oil fields in the case of [14]). For instance, Table III shows that 32:1 compression (1.75 bpp) with a complexity of  $Q = 64\alpha$  yields an average RMSE per band of 2.26, or a little over 2 bits of error. The closest comparable value for non-KLT coding in [14] is a MSE of 40.02 at 2.51 bpp; dividing the MSE by 7 and taking a square root gives an average RMSE per band of 2.39, which is a slightly greater error at a higher bit rate than our result. With interband KLT coding, [14] reports a MSE of 25.11, or an average RMSE of 1.89, at 1.55 bpp; this is a lower distortion at a lower bit rate than our result. We are currently working on incorporating interband KLT coding with the WVQ compression method.

The motivation for our investigation of TM data compression is the need to store and process large amounts of data for postprocessing applications. Using the software package SPECTRUM [15], developed by Los Alamos National Laboratory and the University of New Mexico, we are able to use a desktop workstation running Unix and X-windows to analyze and categorize multispectral data that has been clustered into 256 clusters using a variant of the  $k$ -means algorithm. SPECTRUM can manipulate the color map for the computer display using any transformation of the clustered data, and can display cluster position as a two-dimensional scatter plot. Using these features, users are able to categorize data by selecting areas with a known type of land cover, causing all associated pixels in the image to be given the same pseudocolor representation. Of great interest to us is the robustness of SPECTRUM

data clustering when applied to data that has first been compressed by the WVQ algorithm. While visual quality of pseudocolor visualizations remains good after compression by as much as 32:1, it remains to be determined how much quantization distortion SPECTRUM can tolerate for tasks like Level 1 Land Use Categorization. We are attempting to establish quantitative distortion criteria based on the analysis of classification error presented in [15], which is based on computing levels of confidence for classifications done on clustered data.

#### REFERENCES

- [1] J. N. Bradley, T. G. Stockham, Jr., and V. J. Mathews, "An optimal design procedure for intraband vector quantized subband coding," Tech. Rep. LA-UR-90-4372, Los Alamos National Lab., 1990. Submitted to *IEEE Trans. Commun.*
- [2] J. N. Bradley and C. M. Brislawn, "Image compression by vector quantization of multiresolution decompositions," *Physica D*, vol. 60, pp. 245–258, 1992.
- [3] A. Cohen, I. C. Daubechies, and J.-C. Feauveau, "Biorthogonal bases of compactly supported wavelets," Tech. Rep. 11217-900529-07 TM, AT&T Bell Labs, Murray Hill, NJ, May 1990.
- [4] I. C. Daubechies, *Ten Lectures on Wavelets*. No. 61 in CBMS-NSF Regional Conf. Series in Appl. Math., (Univ. of Lowell, Lowell, MA, June 1990), Philadelphia, PA: Soc. Indust. Appl. Math., 1992.
- [5] C. M. Brislawn, "Classification of symmetric wavelet transforms," Tech. Rep. LA-UR-92-2823, Los Alamos National Lab., Aug. 1992.
- [6] J. N. Bradley, C. M. Brislawn, and V. Faber, "Reflected boundary conditions for multirate filter banks," in *Proc. 1992 IEEE-SP Int. Symp. on Time-Freq. and Time-Scale Analysis*, (Victoria, B.C.), IEEE Signal Processing Soc., Oct. 1992.
- [7] A. Gersho, "Asymptotically optimal block quantization," *IEEE Trans. Info. Theory*, vol. 25, pp. 373–380, July 1979.
- [8] Y. Linde, A. Buzo, and R. M. Gray, "An algorithm for vector quantizer design," *IEEE Trans. Commun.*, vol. 28, pp. 84–95, Jan. 1980.
- [9] A. Gersho and R. M. Gray, *Vector Quantization and Signal Compression*. No. 159 in Int'l. Series in Engineering & Computer Science, Norwell, MA: Kluwer Academic Publishers, 1992.
- [10] J. N. Bradley and C. M. Brislawn, "Wavelet transform-vector quantization compression of supercomputer ocean models," in *Proc. 1993 Data Compress. Conf.*, (Snowbird, UT), IEEE Computer Soc., Mar. 1993. To appear.
- [11] J. N. Bradley and C. M. Brislawn, "Compression of fingerprint data using the wavelet vector quantization image compression algorithm," Tech. Rep. LA-UR-92-1507, Los Alamos National Lab., Apr. 1992. Progress report to the FBI.

- [12] J. K. Dukowicz, R. D. Smith, and R. C. Malone, "A reformulation and implementation of the Bryan-Cox-Semtner ocean model on the Connection Machine," Tech. Rep. LA-UR-91-2864, Los Alamos National Lab., 1991. Submitted to *J. Phys. Oceanogr.*
- [13] R. D. Smith, J. K. Dukowicz, and R. C. Malone, "Parallel ocean general circulation modeling," Tech. Rep. LA-UR-92-200, Los Alamos National Lab., 1992.
- [14] B. R. Epstein, R. Hingorani, J. M. Shapiro, and M. Czigler, "Multispectral KLT-wavelet data compression for landsat thematic mapper images," in *Proc. 1992 Data Compress. Conf.*, (Snowbird, UT), pp. 200-205, IEEE Computer Soc., Mar. 1992.
- [15] P. M. Kelly and J. M. White, "Preprocessing remotely sensed data for efficient analysis and classification," in *Proc. Conf. Knowledge-Based Systems Aerospace Industry*, vol. 1963 of *Proc. SPIE*, (Orlando, FL), Soc. Photo-Opt. Instrument. Engineers, Apr. 1993. To appear.

# Wavelet Encoding and Variable Resolution Progressive Transmission

Ronald P. Blanford  
TRW Systems Integration Group  
One Space Park, Bldg R2/2162  
Redondo Beach, CA 90278

## Abstract

Progressive transmission is a method of transmitting and displaying imagery in stages of successively improving quality. The subsampled lowpass image representations generated by a wavelet transformation suit this purpose well, but for best results the order of presentation is critical. Candidate data for transmission are best selected using dynamic prioritization criteria generated from image contents and viewer guidance. We show that wavelets are not only suitable but superior when used to encode data for progressive transmission at non-uniform resolutions. This application does not preclude additional compression using quantization of highpass coefficients, which to the contrary results in superior image approximations at low data rates.

## 1 Background

Progressive transmission is a method of encoding and transmitting imagery in such a way that gross features are able to be displayed first and subsequently refined to higher and higher resolution. Among the many possible encoding techniques are multiresolution pyramids, discrete cosine transforms, vector quantization, and wavelet transforms. Tzou [6] provides a comprehensive review of proposed techniques for progressive transmission.

The order in which the image data is selected, transmitted, and presented to the user may be dynamically prioritized as a function of both image content and immediate user interest. This typically results in a display which has a non-uniform resolution. Regions containing visually or operationally significant information may be rendered at a much higher resolution, with refinement deferred for areas of uniform intensity or lesser importance. Dreizen [3] proposed one such implementation in which the transmitter identified significant regions and communicated this information to the receiver in addition to the image data. Blanford [2] observed that for a large variety of images this overhead was unnecessary because the receiver could make a reasonable guess at the location of significant regions from image information already transmitted and displayed.

Recent results of compression using wavelet encoding have been shown to provide efficient bit rate reduction while maintaining quite acceptable image quality. The multiresolution nature of the wavelet transform, described by Mallat [5], and its computational efficiency make it a good candidate for fine-grained progressive transmission as well. Antonini *et al.* [1], for example, present

a coarse-grained example of progressive transmission using wavelets by the simple expedient of displaying each lowpass approximation as it is generated during the course of decoding.

In this paper we show that wavelets are not only suitable but superior when used to encode data for fine-grained progressive transmission at non-uniform resolutions. We first describe the approach, then discuss issues and problems in the incorporation of wavelet encoding. We present results showing a marked improvement in the approximations generated for equivalent amounts of data transmitted. Finally we show that the compressible nature of wavelets is not lost in this application; to the contrary, compression by quantization of highpass coefficients results in superior image approximations at low data rates.

## 2 Approach

In a prior publication [2] we presented arguments which led to the conclusion that, in the spatial domain, the low-resolution image approximation which minimizes the mean square error consists of a collection of disjoint regions each of which is painted with the average value of the pixels subsumed. For ease of computation and representation, these regions are restricted to representing nodes in a quadtree constructed by iteratively averaging groups of four pixels. An image approximation, therefore, corresponds to an arbitrary cut through the quadtree, with the minimal approximation being the global average represented by the single node at the apex. The progressive transmission of Antonini *et al.* [1], for example, can be characterized as displaying a set of horizontal cuts corresponding to uniform levels of resolution.

But the cuts need be neither horizontal nor planar. In actuality, the process of transmission may be envisioned as a walk through this quadtree. At each step in the traversal, an unvisited node is selected and expanded by transmitting the information required to generate its children. At the receiving end, the current approximation is transformed into its successor by using the new information to generate and paint the values of the child regions. Thus each approximation differs from its predecessor in that a single region has been replaced by four subregions. The traversal terminates when all regions are leaf nodes one pixel in size.

The question then arises which of the many possible traversals is optimal. For a non-interactive transmission the goal might be to minimize the mean square error, in which case a greedy algorithm can be applied at each step to select for refinement the region with maximum error. In an interactive situation, if the viewer has indicated a particular point of interest in the image then the traversal might select the non-leaf region nearest that point. If a particular feature is of interest, then the region could be selected which responds most strongly to a feature detection algorithm. Or the selection might be based on a combination of several criteria.

In this discussion we will model a non-interactive session with the goal of minimizing the error represented in the approximation. The region with the greatest error is that whose product of pixel variance and area is maximized. The receiver knows the area but not the pixel variance of the regions in its approximation. The transmitter knows both and could send a region identifier along with the information needed to refine the region, but this overhead is generally unacceptable in the low-bandwidth situations where progressive transmission is most useful. It turns out, however, that for a wide variety of images a good predictor of the pixel variance within a region is the variance between the region and those which neighbor it in the current approximation. The receiver and transmitter can independently perform this computation to select the region and only the corresponding image

information need be transmitted.

### 3 Image encoding

In the previous work, we introduced the additional constraints that the encoding method used to build the quadtree be lossless and introduce no storage overhead. These constraints led to the selection of the *comp/diff* encoding scheme first proposed by Knowlton [4]. In this paper we will relax first the overhead constraint and then the lossless constraint and show by comparison how the resulting approximations fare.

Knowlton’s *comp/diff* encoding applied to two pixel values produces a composite value which approximates the average and a differentiator value which approximates the difference. Each of these values requires precisely the same number of bits as the original pixels, and so requires no storage overhead. The same encoding function applied to the composite and differentiator returns the original pixels, so the procedure is lossless. All is not rosy, however, as the encoded composite value may differ significantly from the true average. The error is exacerbated if quantization is attempted in an effort at compression.

The Knowlton encoding is a non-linear function but resembles a wavelet transform with two taps: low frequency information is captured in the composite while high-frequency information resides in the differentiator. The current experiment replaces that encoding with a wavelet of eight taps. The highpass coefficients generated by the wavelet transform require roughly the same number of bits for representation as the original input, but the lowpass coefficients typically require one additional bit. Thus the encoding is not without storage overhead, which empirical evidence shows can be as high as twenty percent. Most of the additional bits are used to represent the top levels of the quadtree which are transmitted first, so the number of coefficients transmitted will be fewer than with the same amount of data using the Knowlton encoding. Our hope is that the quality of the coefficients will more than compensate for the lesser number.

Constructing the encoded quadtree presents no undue difficulties. We choose to treat the image in a toroidal manner, wrapping left to right and top to bottom, so that it will not be necessary to add extra rows and columns of padding at the lower resolutions. We use a separable wavelet function having one-dimensional 8-element low-pass and high-pass kernels.

Because of the larger basis we can no longer construct the entire quadtree, but must stop at the second (4x4) level below the apex. These sixteen lowpass values are transmitted and displayed as the initial approximation. For each region we compute a refinement priority which is the product of the region area and its external variance. The area is just the number of pixels the region represents. The external variance of region  $R$  is the mean square difference between the region value  $p_R$  and the values of regions found in parts of the displayed approximation immediately adjacent to the region. Let us compute an estimate  $v_R^s$  of the variance using those neighbors, where  $w_N$  is the length of the side shared with neighbor  $N$  and  $p_N$  is its displayed value.

$$v_R^s = \frac{1}{4w_R} \sum_{N \in \text{neighbors}(R)} w_N (p_R - p_N)^2 \quad (1)$$

After assigning priorities to all initial regions, we enter the refinement phase.

## 4 Image refinement

The refinement proceeds in three steps, iterated repeatedly until transmission is complete.

1. Select the maximum priority region for refinement.
2. Transmit the encoded data needed to produce the four subregions.
3. Compute priorities for the new subregions and their immediate neighbors.

The selection of the maximum priority region is simple at the beginning of the transmission, but with a brute-force approach would quickly grow intransigent as the number of regions multiplies to a significant fraction of the image size. We deal with this problem by creating and maintaining a priority heap which holds all unvisited regions. The effort to insert and update priorities in a heap with  $N$  entries is of order  $O(\log N)$ , which renders the problem manageable. The next region to be selected is always at the top of the heap.

With the Knowlton encoding, the region value together with three additional differentiators was all that was needed to compute the values of its four children. The broader basis of the wavelet encoding necessitates that, if region  $R$  is to be refined, then not only must its value be present but also the lowpass values of neighbors in a  $5 \times 5$  area surrounding it. If a neighbor is found whose value has not yet been computed, its parent is selected for refinement regardless of priority and its  $5 \times 5$  neighborhood examined for data availability. Because the initial regions all satisfy this neighborhood criterion, the procedure must eventually succeed in selecting a region.

The second step is to identify and transmit the additional information needed to produce the four children. Just as with the Knowlton encoding, we require three highpass coefficients for each lowpass one. The difference is that we now require 25 times as many, and that some may have already been transmitted for use with other regions. The bookkeeping required to determine which coefficients remain to be transmitted is painstaking but not unduly taxing. Once all required coefficients have been provided, the effort to decode the coefficients and produce the subregion values is trivial. These four values then replace the original region value in the displayed approximation.

Finally priorities must be computed for the newly created regions. Also, since these new region values contribute to the priority computations of their immediate neighbors, their priorities must be invalidated and recomputed as well. The replaced region is removed from the priority heap and all new priorities inserted or updated. This process adjusts the heap so that the new maximum priority moves to the top, in preparation for the next iteration.

## 5 Results of wavelet encoding

Figure 1 shows an aerial view of a portion of Moffet Field, California. The original image contains 8-bit data, 512 pixels on an edge. There are many small features as well as sharp edges between foreground and background, which make it a rather difficult image to compress effectively. Figure 2 gives a graphic comparison of the approximation error as the transmission progresses. The horizontal axis measures the number of bits transmitted as a percentage of the original image size. Knowlton encoding is shown in dark gray. Wavelet encoding is shown in black. The graph clearly

shows that wavelet encoding provides approximations with half the error over most of the transmission, a significant improvement. The light gray curve is the result of transmitting the wavelet coefficients at a uniform resolution: a breadth-first traversal of the quadtree. The knees in the curve correspond to resolution changes. The variable resolution approach is clearly superior.

It is difficult to appreciate the actual impact of progressive transmission in a static presentation such as this. Figures 3 and 4 show snapshots of the display as it would appear when 2% and 5% of the data has been transmitted, using the wavelet encoding. Figures 5 and 6 show the corresponding displays using the Knowlton encoding. The Knowlton encoding provides higher contrasts and sharper edges where they are found, but the wavelet encoding provides a more balanced development.

Figures 7 and 8 show results of transmission error for an aerial view of the Los Angeles airport. The image itself is hazy with few contrasts so it should exhibit lower overall levels of error, as the graph bears out. The black curve for variable resolution wavelet encoding again shows half the error of the dark gray Knowlton encoding. The light gray curve for breadth-first transmission is close to that for variable resolution, as one would expect when few features stand out.

Figures 9 and 10 show results for yet another aerial view of an airport, this one a Spot satellite image of Beirut, Lebanon. Though the overall intensity and feature distribution are different, this image exhibits approximation error similar to the Los Angeles image.

## 6 Results of coefficient quantization

One reason wavelets provide a good basis for image compression is that the result degrades gracefully under quantization, even at extreme levels. In order to verify that this characteristic had not been lost when used with progressive transmission, a simple quantization scheme was applied to the wavelet coefficients and the impact on the approximation error observed. Briefly, this scheme divides all highpass coefficients in a given level of the quadtree by the same value. This value is greatest at the lowest level, the leaf nodes, and is reduced by a factor of two at each level above.

Figure 11 shows the results when applied to the Moffet Field image. The thick black curve on the right represents the uncompressed wavelet transmission, just as in Figure 2, though shown at an expanded scale. The second curve from the right represents a quantization factor of 4 at the base level of the quadtree, of 2 at the level above the base, and the remainder of the coefficients left unquantized. The base-level quantization factor is doubled for each successive curve, so that the leftmost curve begins with a quantization factor of 64 at the base level of the quadtree, reduced successively to a divisor of 2 at level four, with only the original lowpass coefficients remaining unscathed.

While the higher levels of quantization introduce significant error in the later stages of the transmission, the initial portions critical for early identification of features show only improvement. More significantly, the smooth degradation indicates that efforts toward designing more sophisticated quantization schemes would not go unrewarded.

## 7 Conclusions

In this paper we have shown that wavelets are not only suitable but superior when used to encode data for fine-grained progressive transmission at non-uniform resolutions. The results show a marked improvement in the approximations generated for equivalent amounts of data transmitted when wavelet encoding is used in place of Knowlton encoding. Finally we have shown that the compressible nature of wavelets is not lost in this application; to the contrary, compression by quantization of highpass coefficients results in superior image approximations at low data rates.

## References

- [1] M. Antonini, M. Barlaud, P. Mathiew, and I. Daubechies. Image coding using wavelet transform. *IEEE Transactions on Image Processing*, 1:205–220, 1992.
- [2] Ronald P. Blanford. Progressive refinement using local variance estimators. *IEEE Transactions on Communications*, to appear.
- [3] Howard M. Dreizen. Content-driven progressive transmission of grey-scale images. *IEEE Transactions on Communications*, COM-35(3):289–296, March 1987.
- [4] Ken Knowlton. Progressive transmission of grey-scale and binary pictures by simple, efficient, and loss-less encoding schemes. *Proceedings of the IEEE*, 68(7):885–896, July 1980.
- [5] S. Mallat. A theory of multiresolution signal decomposition: The wavelet representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11:674–693, 1989.
- [6] Kou-Hu Tzou. Progressive image transmission: A review and comparison of techniques. *Optical Engineering*, 26(7):581–589, July 1987.

Moffet Field data:

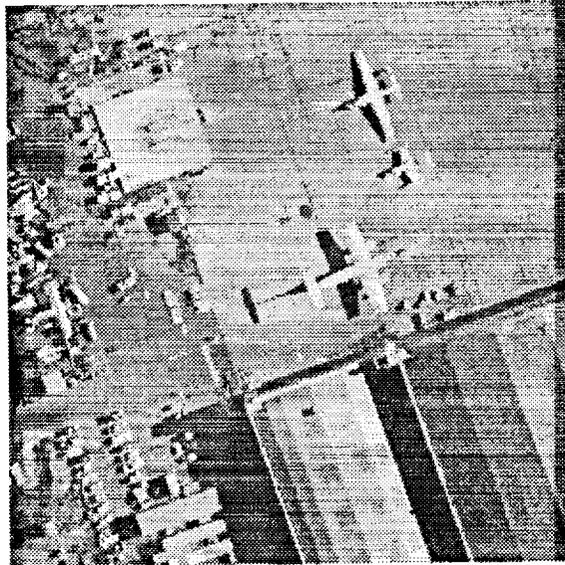


Figure 1: Moffet Field original image

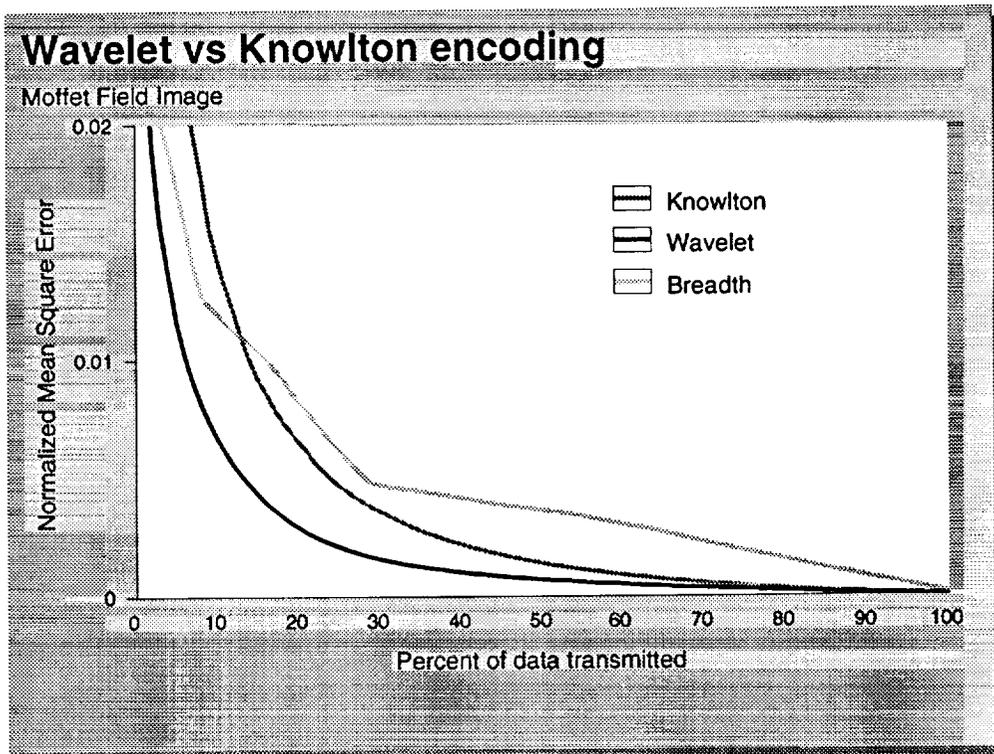


Figure 2: Comparison of encoding error

**Moffet Field snapshots:**



Figure 3: Wavelet encoding, 2% of data

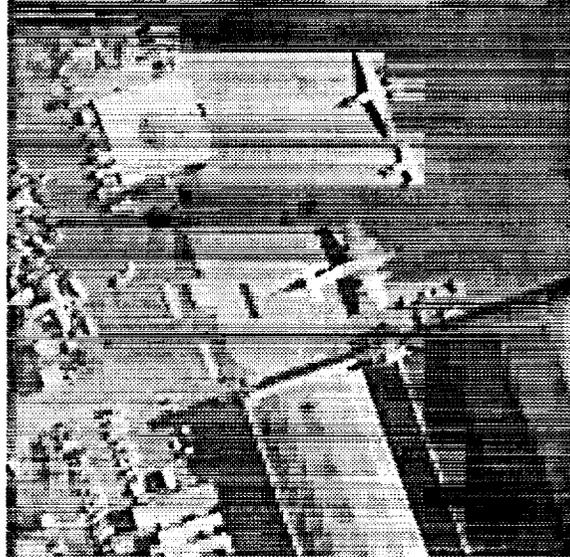


Figure 4: Wavelet encoding, 5% of data

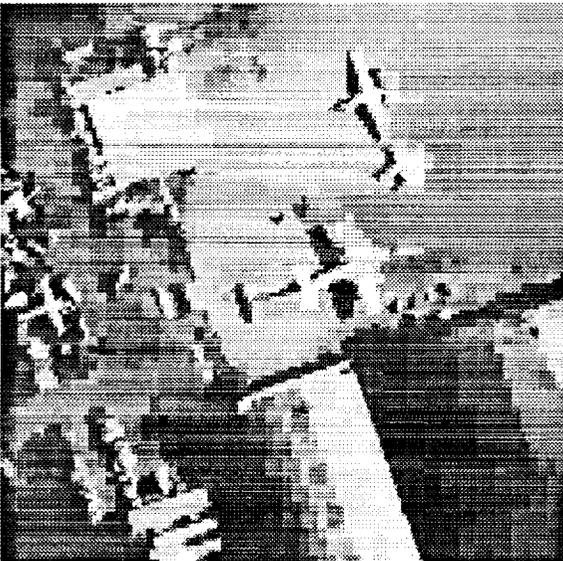


Figure 5: Knowlton encoding, 2% of data

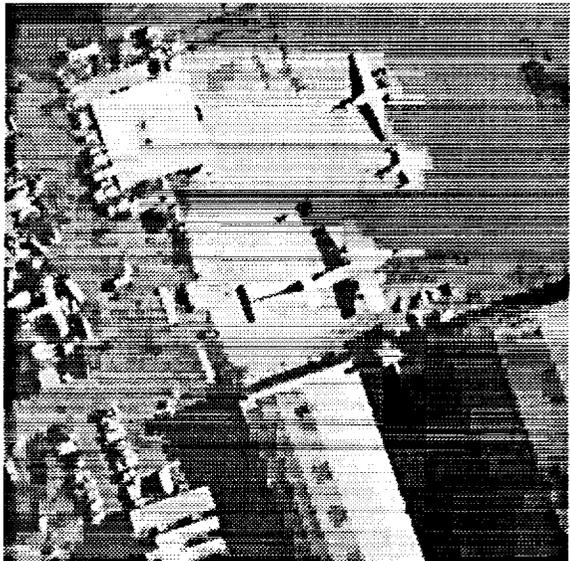


Figure 6: Knowlton encoding, 5% of data

Los Angeles data:

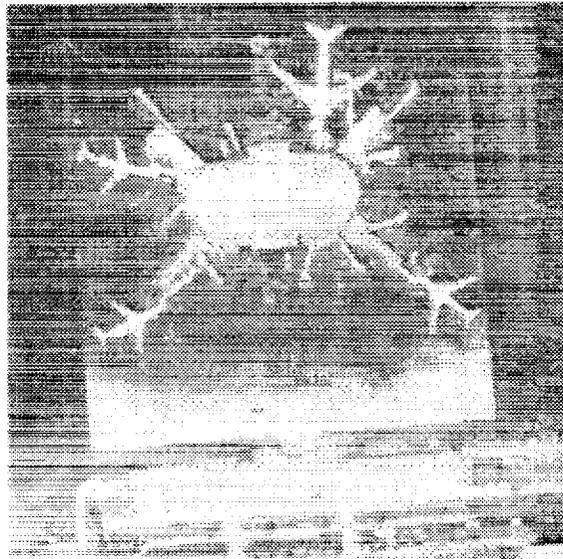


Figure 7: Los Angeles Airport original image

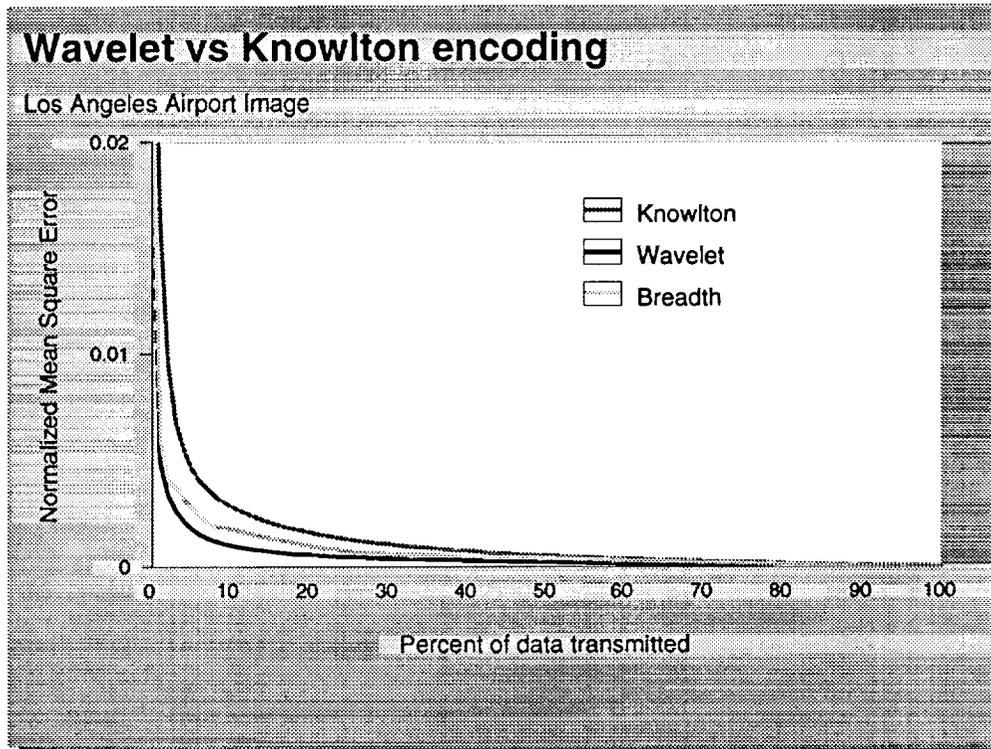


Figure 8: Comparison of encoding error

Beirut Airport data:

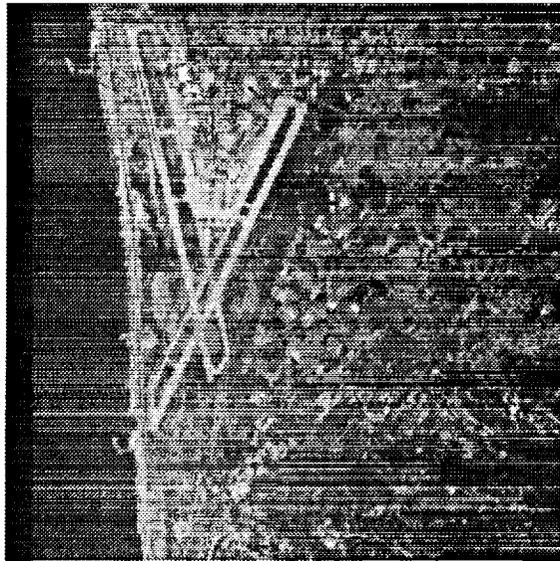


Figure 9: Beirut Airport original image

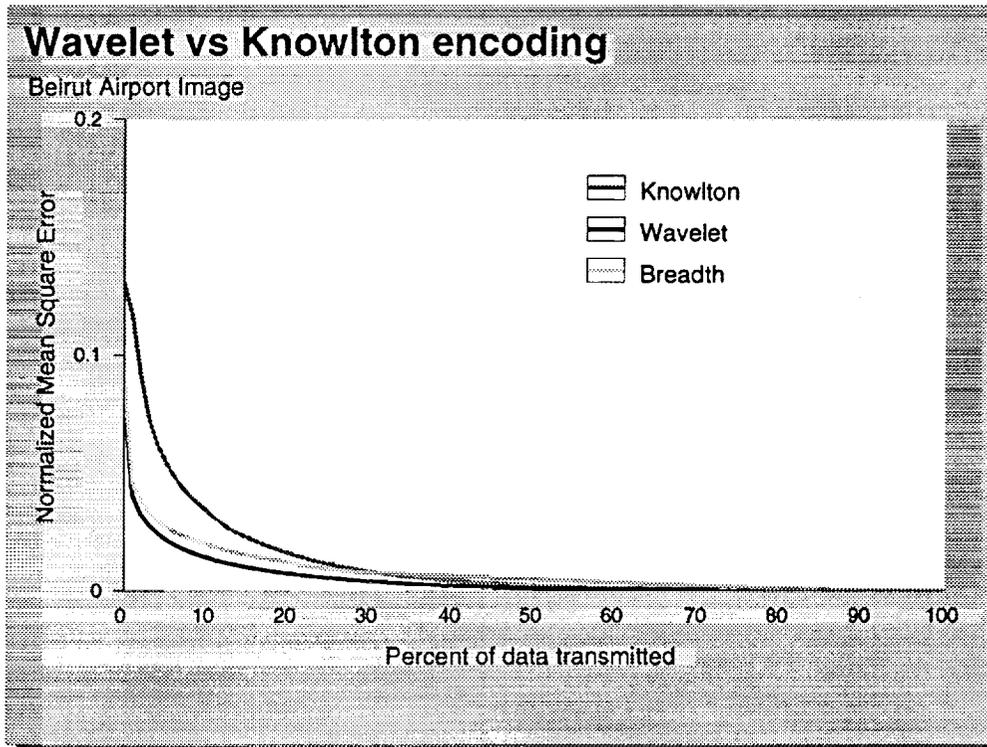


Figure 10: Comparison of encoding error

Quantization data:

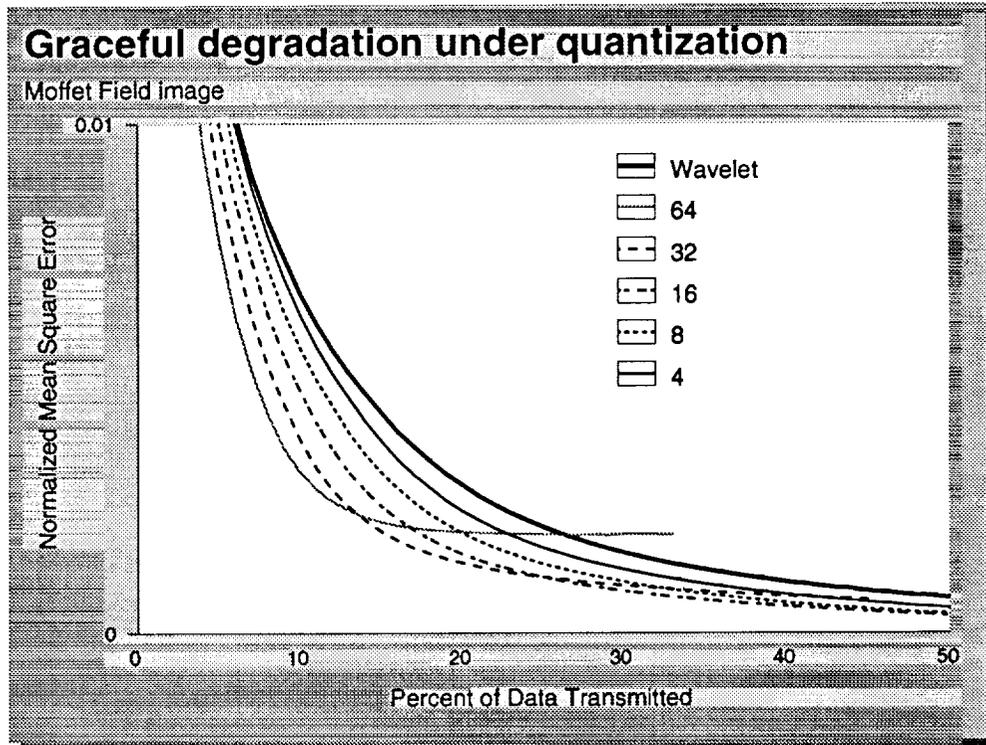


Figure 11: Comparison of quantization error



## FAST IMAGE DECOMPRESSION FOR TELEBROWSING OF IMAGES

Shaou-Gang Miaou and Julius T. Tou

Department of Electrical Engineering  
University of Florida, Gainesville, FL 32611

**Abstract.** Progressive image transmission (PIT) is often used to reduce the transmission time of an image telebrowsing system. A side effect of the PIT is the increase of computational complexity at the viewer's site. This effect is more serious in transform domain techniques than in other techniques. Recent attempts to reduce the side effect are futile as they create another side effect, namely, the discontinuous and unpleasant image build-up. Based on a practical assumption that image blocks to be inverse transformed are generally sparse, this paper presents a method to minimize both side effects simultaneously.

### 1. Introduction

One important evaluation criterion for a telebrowsing system is the response time which is the time elapsed from the moment a retrieval request is issued until the desired information is actually displayed on the monitor [1]. The response time can roughly be divided into three major parts. The first part is the searching time for the system to locate the desired information. The second part is the transmission time to send the information through a channel. The third part is the display time for the information to be displayed on the monitor. The early studies of the telebrowsing systems were concentrated on the efficient retrieval of pure text information [2,3]. In this case, the searching time is the only major concern. However, for modern telebrowsing systems where multimedia information, including text, audio, image, and video, is considered, the transmission time and the display time become a significant part of the response time because of huge amount of data involved in still images and video (a sequence of images).

To reduce the transmission time of an image telebrowsing system, a well known scheme called progressive image transmission (PIT) is often used. PIT allows an approximate reconstruction of an image whose fidelity is built up gradually until the viewer decides either to abort the transmission sequence or to allow further reconstruction. This scheme increases the effective compression ratio because usually only a small part of the compressed data needs to be sent for browsing purpose.

With PIT techniques, the transmission time can be greatly reduced. However, it also creates a side effect, that is, it increases the processing time at the viewer's site because an inverse PIT process is required. Since the major task of the inverse PIT process is the image decompression given part of the compressed data, the research is aimed at the development of fast image decompression schemes for the inverse PIT process.

The rest of the paper is organized as follows. First, the PIT schemes and their computational complexities are briefly addressed. Then, the drawbacks of recent attempts to reduce the computational complexities are discussed. Next, the demonstration of a new approach is given. Finally, a performance comparison between the new approach and the recent ones is made.

## 2. PIT Schemes and Their Computational Complexities

There are many PIT schemes. According to Tzou's classification, they are divided into three major categories, namely, spatial domain, transform domain, and pyramid-structured, based on where the progression takes place [4]. Each category can be further divided into several classes of techniques. The classification is shown in Figure 1. Note that not all of the PIT schemes will produce a considerable amount of computational overhead in the inverse PIT process. For instance, the spatial domain schemes only require a very low computational effort in the inverse PIT process. In pyramid-structured PIT schemes, only successively filtered pyramid techniques require high computational complexity in the inverse PIT process. Even for the successively filtered pyramid techniques, however, the complexity to process the first few levels of a pyramid from the top remains low. From a practical point of view, the processing of the first few levels of the pyramid may suffice the purpose of image browsing. On the other hand, transform domain techniques usually take considerable amount of computation in the inverse PIT process, since the inverse transforms have to be carried out with about the same computational effort for every stage of image reconstruction.

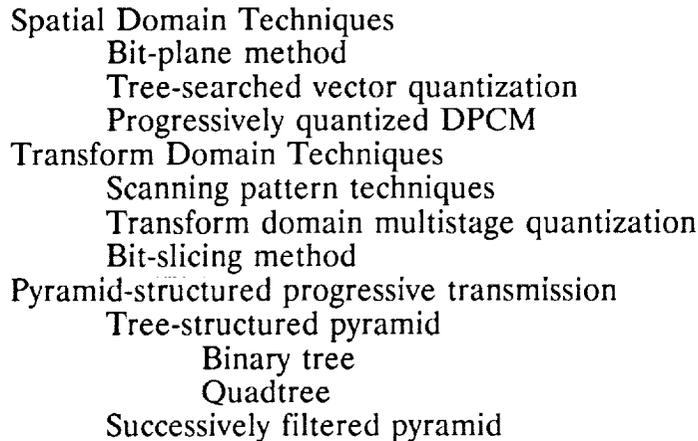


Figure 1. Tzou's classification of PIT schemes

In transform domain PIT schemes, the transform coefficients are first quantized and then divided into segments. Only one segment of quantized coefficients is sent for one stage of image reconstruction. The only differences among all transform domain techniques are the ways to determine the segments and the order in which they are sent. One common feature among them is that the transform coefficients are only "partially" encoded where in a non-PIT or sequential scheme they are said to be "fully" encoded.

For a transform domain non-PIT scheme, one  $M \times N$  inverse transform for an  $M \times N$  image block is needed. However, for a transform domain PIT scheme,  $r$  times of  $M \times N$  inverse transform are needed for the image block, where  $r$  is the number of stages of image reconstruction. The lower bound for  $r$  is 1 but its upper bound depends on the image, the viewer, and the PIT scheme. Therefore, the computation load for inverse transform is  $r$  times heavier in PIT schemes than in non-PIT schemes.

One transform domain scheme using discrete cosine transform (DCT) receives great attention, since the DCT has the energy packing capabilities and also approaches the

statistically optimal transform (i.e. Karhunen-Loeve transform) in decorrelating a signal governed by a Markov process [5]. In addition, it is part of the recently approved JPEG standard [6,7]. The JPEG standard has brought a tremendous impact on the image-coded related industry. However, as far as implementation of the standard is concerned, the standard provides only a guideline. How to implement the standard efficiently for certain application still relies on the ingenuity of designers. For example, JPEG has chosen to specify neither a unique forward DCT (FDCT) algorithm or a unique inverse DCT (IDCT) in its recommendation. This is because research in fast DCT algorithms is ongoing and no single algorithm is optimal for all implementations [7]. For the application of inverse PIT, we will show that traditional fast two dimensional (2-D) IDCT algorithms can be accelerated to reduce the processing time at the viewer's site.

### 3. Previous Approaches and Their Drawbacks

To relieve the computation burden of IDCTs in inverse PIT, the following approaches have been used.

Approach 1: Use traditional fast algorithms for IDCT. The computational complexity is reduced from  $O(N^4)$  by the definition of IDCT to  $O(N^2 \log_2 N)$  by traditional fast algorithms, where  $N \times N$  is the block size. There are many fast algorithms available for IDCT. For an  $8 \times 8$  IDCT, one of the best algorithms reported so far takes 96 multiplications and 466 additions [8].

Approach 2: Use a fast progressive reconstruction method. It is a combination of a special scheme and the use of approach 1. This approach was first proposed by Takikawa to perform fast progressive reconstruction for discrete Fourier transformed and Walsh-Hadamard transformed images [9]. Later, Miran and Rao followed the similar derivation by Takikawa and developed a fast progressive reconstruction for DCT images [10]. The basic idea of approach 2 is to decompose the  $N \times N$  transformed block into  $\log_2 N + 1$  sparse matrices, each of which can be inverse transformed by  $1 \times 1$ ,  $2 \times 2$ ,  $4 \times 4$ , ..., and  $N \times N$  fast inverse transform algorithms.

Approach 2 has some advantages over approach 1. First, the computational complexity is lower. For example, consider a 4-stage image reconstruction and an  $8 \times 8$  image block. Approach 1 takes four  $8 \times 8$  IDCTs while approach 2 requires only one  $1 \times 1$  IDCT, one  $2 \times 2$  IDCT, one  $4 \times 4$  IDCT, and one  $8 \times 8$  IDCT. The computational saving is obvious. Secondly, the delay time is reduced. The delay time is the time to wait for all the elements in a transformed block before an inverse fast transform can be performed.

However, approach 2 has a serious problem, that is, it has a poor and discontinuous image build up. The reason is that the order in which the sparse matrices are formed and sent is not in the order of visual significance. In general, a DCT coefficient with higher variance (or energy) tends to be more visually significant than that with lower variance. It is well known that the DCT coefficient variances are highly correlated along the zig-zag scan [11]. Approach 2 has a fixed transmission pattern that does not even close to the zig-zag scan. This problem has been confirmed experimentally by Miran and Rao [10]. They ascribed the drawback to not having low frequency terms immediately adjacent to DC components in the intermediate stages of reconstruction. Another drawback of approach 2 is that it still

requires all elements of the sparse matrices to start computing the inverse transform. Thus, the delay time is reduced but not eliminated.

One more drawback for both approaches 1 and 2 is the computational redundancy of traditional fast algorithms in inverse PIT. If IDCT is used in image decompression, its input block contains only a few nonzero coefficients. In addition, if a PIT scheme is used, the input matrix to IDCT contains even fewer nonzero elements. To visualize the redundancy, consider the signal flowgraph of a fast IDCT algorithm. Since a zero presented at an input node contributes nothing to the output, the paths between a zero input node and output nodes are trivial or redundant.

To get a picture on how many spatial frequencies retained on the average after the quantization, many 512 x 512 8-bit greyscale images and RGB components of color images were tested. In the test, JPEG's coding scheme, including a recommended quantization table, was used. Part of the test result is presented in column 3 of Table 1. It is shown, even for a very busy image such as baboon image, no more than a quarter of quantized coefficients are nonzero. Even with so many spatial frequencies set to zero, the decompressed images and their originals are perceptually indistinguishable. Next, consider the case when a small visible image degradation is allowed. To produce a small image degradation, the same test was repeated except that the round-off operation in JPEG's scheme was replaced by truncation. The new numbers are shown in column 4. The decompressed images have only minor degradation, for it does not diminish our capability to recognize meaningful objects in the images. In fact, the image quality is good enough to be the last stage of PIT. The test shows that the number of nonzero quantized DCT coefficients decreases sharply at the minor expense of image quality. In the inverse PIT process, the average number of nonzero elements in an 8 x 8 matrix does not need to be higher than that in column 4.

Table 1. Average # of Nonzero Quantized DCT Coefficients in an 8 x 8 Block

Images	Image Activity	Round-off	Truncation
Lena	Low	6.13	4.01
Boat	Medium	9.20	6.00
Baboon	High	15.50	9.80

How much of the matrix must be zero for it to be considered sparse depends on the applications. Generally, a matrix is called sparse if there is an advantage in exploiting its zeros [12]. It is well known that exploiting the sparsity can lead to enormous computational savings in many applications such as solving simultaneous equations with Gaussian elimination method. Inspired by this fact, it is curious to see if the sparsity of the input matrix can also be exploited to compute IDCT efficiently in the environment of inverse PIT. Since the characteristic of an input image block to IDCT is generally not considered in traditional algorithms, a nonconventional approach must be adopted to exploit the sparsity of the input matrix. The proposed approach will be presented in the following manner. First, we describe the goal to be accomplished by the approach. Then, the rationale of the approach is discussed. Next, based on the rationale, two methods are presented -- one is too slow to be useful, the other is its fast version. The fast version is shown to be good enough for the practical use.

#### 4. The Proposed Approach

In the inverse PIT process, computation burden of IDCT and computation redundancy associated with traditional algorithms are two major problems. The inherent drawbacks in Takikawa's or Miran and Rao's approach present another problem in the inverse PIT process. In view of all these problems, our approach should meet the following goals. First, it must be fast and efficient. Second, it must allow a scanning pattern that can conform to the visual significance. Finally, it must have practically no delay time.

For the ease of discussion, several terms are defined first. A target matrix is an image block consisting of the quantized DCT coefficients that are partially encoded for PIT. Performing an IDCT on a target matrix results in a matrix called goal matrix. The result of processing one nonzero element in the target matrix is called the partial contribution to the goal matrix. Throughout this paper, the partial contribution is treated as a matrix or all its elements depending on the context.

Based on the definition of 2-D IDCT, only nonzero elements in the target matrix can contribute to the goal matrix. In fact, the value of each nonzero element can affect the values of all elements in the goal matrix. The idea of our approach is to completely ignore the zero elements in the target matrix and process each nonzero element separately and efficiently. The goal matrix is then updated periodically by adding the partial contribution. Therefore, the computation of IDCT is divided into two tasks, i.e., the computation of partial contribution and the update of the goal matrix. The idea adapts particularly well to the scheme where DCT coefficients are run-length coded (such as JPEG's).

The definition of 2-D IDCT is

$$g_{xy} = \frac{2}{\sqrt{MN}} \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} c(u)c(v)F_{uv} \cos\left(\frac{(2x+1)u\pi}{2M}\right) \cos\left(\frac{(2y+1)v\pi}{2N}\right) \quad (1)$$

where  $x=0, \dots, M-1, y=0, \dots, N-1$ , and

$$c(k) = \begin{cases} \frac{1}{\sqrt{2}} & \text{if } k=0, \\ 1 & \text{otherwise.} \end{cases}$$

The coefficient in front of the double summation of equation 1 is only a scale factor which requires essentially no computation (except a register shift operation) in practical applications where  $M=N$  and  $M=4, 8, \text{ or } 16$  are often used. Thus, it is usually neglected when comparing the computational complexity among the fast algorithms of IDCT. By taking the scale factor out, equation 1 becomes

$$f_{xy} = \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} c(u)c(v)F_{uv} \cos\left(\frac{(2x+1)u\pi}{2M}\right) \cos\left(\frac{(2y+1)v\pi}{2N}\right) \quad (2)$$

where

$$f_{xy} = \frac{\sqrt{MN}}{2} g_{xy}$$

If  $[f_{xy}]_{uv}$  is defined as the partial contribution to the goal matrix  $[f_{xy}]$  due to  $F_{uv}$  alone, then

$$[f_{xy}]_{uv} = c(u)c(v)F_{uv}\cos\left(\frac{(2x+1)u\pi}{2M}\right)\cos\left(\frac{(2y+1)v\pi}{2N}\right) \quad (3)$$

The partial contribution can be obtained by the use of definition in equation 3 directly. Assume the values of cosine functions for different combinations of  $x$  and  $u$  are precalculated and stored as a table. The table can also be used as the values of cosine functions for different combinations of  $y$  and  $v$  with  $x$  and  $u$  replaced by  $y$  and  $v$ , respectively. Let  $Q$  be the number of multiplications required to find the partial contribution due to  $F_{uv}$ . Then,  $Q=2MN$  if both  $u$  and  $v$  are not zero,  $Q=3MN$  if  $u=0$  or  $v=0$  but not both, and  $Q=0$  if  $u$  and  $v$  are both zero. For an  $M \times N$  target matrix with  $n$  ( $> 1$ ) nonzero  $F_{uv}$ , where  $n \ll MN$ , the number of multiplications to get the goal matrix is from  $2(n-1)MN$  to  $3nMN$ . With this naive approach, no addition but 128 to 384 multiplications are required if  $M=N=8$  and  $n=2$ . This is not good enough, since an  $8 \times 8$  fast IDCT can take as low as 96 multiplications [8]. Therefore, a better way to compute the partial contribution is needed.

Equations 2 and 3 are equivalent if only one term in the double summation of equation 2 is nonzero. So the traditional fast algorithms for equation 2 can be applied to equation 3 as well. However, the direct use of them to compute the partial contribution is not desirable since they contain high computational redundancy. We found that with a systematic reduction rule for the signal flowgraphs of traditional fast algorithms, a much faster way to compute the partial contribution than the naive approach is possible. The rule is based on the two attributes associated with the partial contribution, which we call the mirror effect and the reducible property.

From equation 3, it can be readily shown that

$$\begin{aligned} [f_{x'y'}]_{uv} &= (-1)^u [f_{M-1-x',y'}]_{uv} \\ [f_{xy'}]_{uv} &= (-1)^v [f_{x,N-1-y'}]_{uv} \\ [f_{x'y'}]_{uv} &= (-1)^{u+v} [f_{M-1-x',N-1-y'}]_{uv} \end{aligned} \quad (4)$$

where  $x'$  and  $y'$  are particular values of  $x$  and  $y$ , respectively. Equation 4 indicates that the partial contribution exhibits high degree of symmetry or mirror effect. Note that only possible sign changes are involved in equation 4 and practically no addition or multiplication is required. The significance of this result is that only a quarter of the partial contribution needs to be determined through additions and multiplications. The rest of them can be determined by simple copy operations and possible sign changes.

The reducible property can be stated as follows. The  $M \times N$  partial contribution due to  $F_{uv}$  is equivalent to that of  $(M/m) \times (N/n)$  partial contribution due to  $F_{u/m,v/n}$ , where  $cd(M,u) = m$ ,  $cd(N,v) = n$ , and  $cd(a,b)$  is a common divisor between non-negative integers  $a$  and  $b$ . This statement can be proved easily by noting that

$${}_{M,N}[f_{xy}]_{u,v} = {}_{M/m,N/n}[f_{xy}]_{u/m,v/n} \quad (5)$$

where  ${}_{M,N}[f_{xy}]_{u,v}$  is the  $[f_{xy}]_{uv}$  defined in equation 3. If  $cd(M,u)=1$  only, it is said to be an irreducible partial contribution in row. Similarly, if  $cd(N,v)=1$  only, it is said to be irreducible in column. If for some  $m>1$  or  $n>1$ , the partial contribution is said to be reducible. Note that the reducible property is separable, i.e, the reduction in row size and column size can be processed separately. The largest reachable reduction for  $F_{uv}$  happens

when  $\gcd(M,u)=m$  and  $\gcd(N,v)=n$ , where  $m>1$  and  $n>1$ , and  $\gcd(a,b)$  is the greatest common divisor between non-negative integers  $a$  and  $b$ . For example, if  $M=N=8$ ,  $u=6$ , and  $v=4$ , the  $8 \times 8$  partial contribution due to  $F_{64}$  is equivalent to  $4 \times 2$  partial contribution due to  $F_{31}$ , since  $\gcd(8,6)$  can be 2 and  $\gcd(8,4)$  can be 4. It is also a maximum reducible case for  $F_{64}$  since  $\gcd(8,6)=2$  and  $\gcd(8,4)=4$ . Note that  $\gcd(a,0) = a$ . Therefore, if  $u$  and  $v$  are both zero,  $M \times N$  partial contribution due to  $F_{uv}$  is simply a  $1 \times 1$  partial contribution due to  $F_{00}$ , which is always  $F_{00}/2$  no matter what the values of  $M$  and  $N$  are.

Combining the mirror effect and the reducible property can lead to a great saving in computation of partial contribution. Consider the following example: We want to compute the  $8 \times 8$  partial contribution due to  $F_{44}$ . Since  $\gcd(8,4)=4$ , it can be reduced to the  $2 \times 2$  partial contribution due to  $F_{11}$ . By using the mirror effect, only  $1 \times 1$  of the  $2 \times 2$  partial contribution needs to be determined explicitly, which is  $f_{00}=F_{44}/2$ . It can then be expanded to  $2 \times 2$  partial contribution by the use of mirror effect:

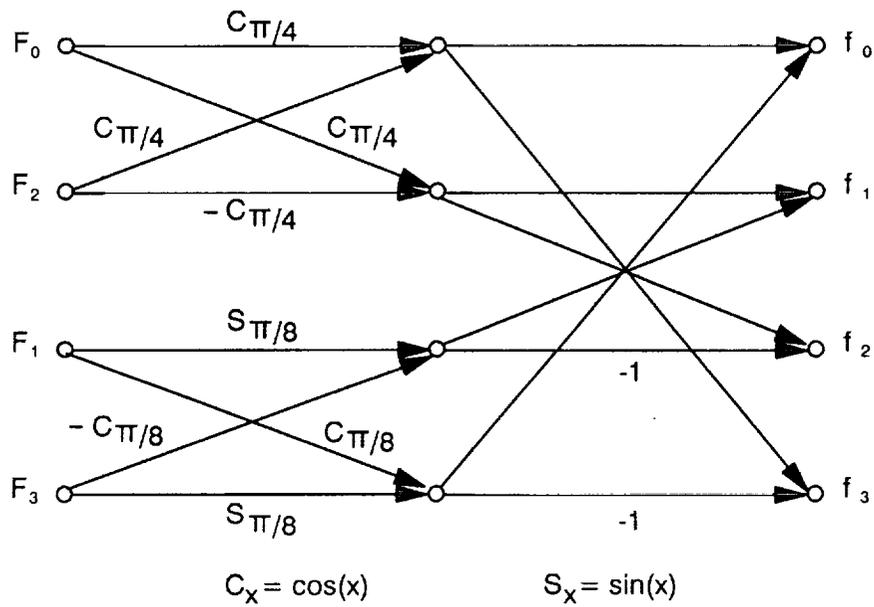
$$\begin{matrix} f_{00} & f_{01} \\ f_{10} & f_{11} \end{matrix}$$

where  $f_{01}=-f_{00}$ ,  $f_{10}=-f_{00}$ , and  $f_{11}=f_{00}$ . Similarly, by using the mirror effect for  $F_{11}$ , we expand the partial contribution to  $4 \times 4$  :

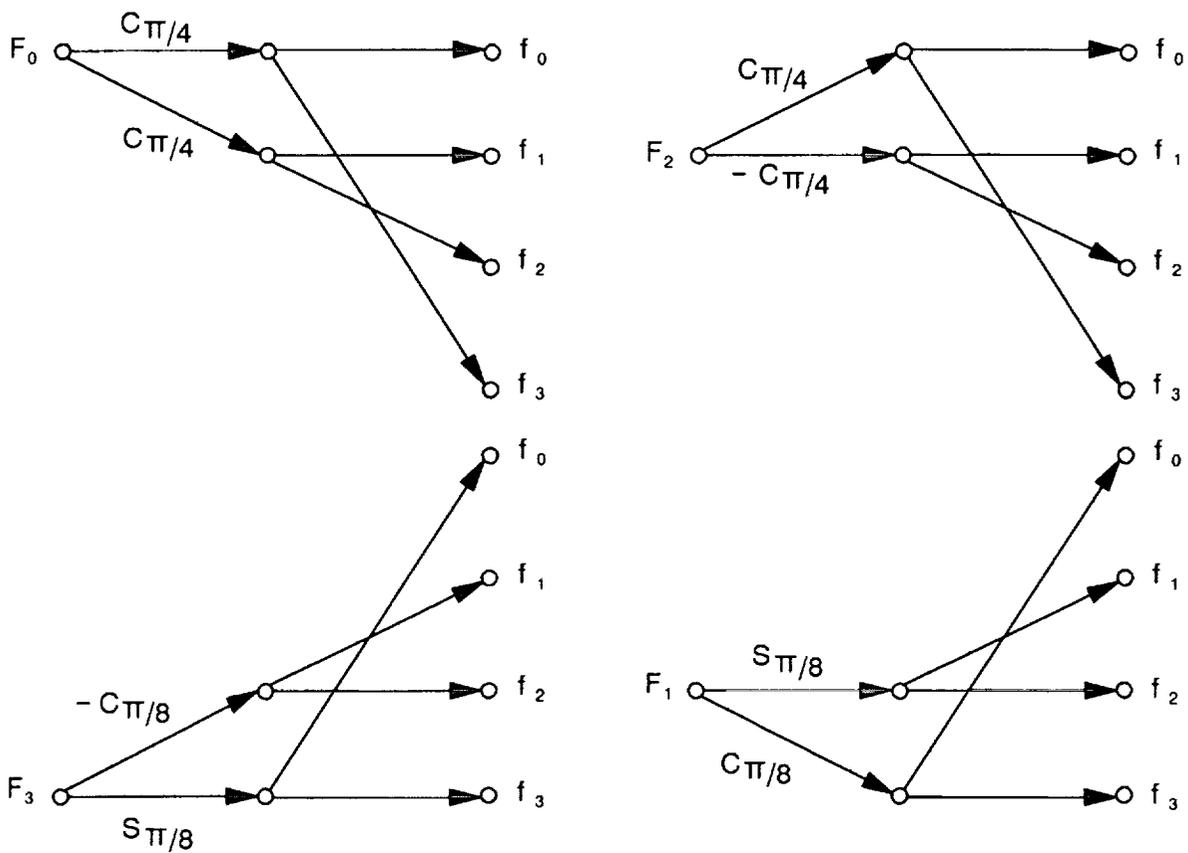
$$\begin{matrix} f_{00} & f_{01} & f_{02} & f_{03} \\ f_{10} & f_{11} & f_{12} & f_{13} \\ f_{20} & f_{21} & f_{22} & f_{23} \\ f_{30} & f_{31} & f_{32} & f_{33} \end{matrix}$$

where  $f_{02}=-f_{01}$ ,  $f_{03}=-f_{00}$ ,  $f_{12}=-f_{11}$ ,  $f_{13}=-f_{10}$ ,  $f_{20}=-f_{10}$ ,  $f_{21}=-f_{11}$ ,  $f_{30}=-f_{00}$ ,  $f_{31}=-f_{01}$ ,  $f_{22}=f_{11}$ ,  $f_{23}=f_{10}$ ,  $f_{32}=f_{01}$ ,  $f_{33}=f_{00}$ . Since  $\gcd(8,4)$  can be 2, the  $8 \times 8$  partial contribution due to  $F_{44}$  is equivalent to the  $4 \times 4$  partial contribution due to  $F_{22}$ . By using the mirror effect for  $F_{22}$ , we can expand the partial contribution of  $4 \times 4$  to the desired result of  $8 \times 8$ . Note that no multiplication is required to determine the 64 elements of partial contribution due to  $F_{44}$ .

The basic principle to reduce the signal flowgraph of a traditional algorithm is by retaining only the nontrivial paths. This concept is demonstrated by an example. Consider the row-column or indirect approach of a fast 2-D IDCT for a  $4 \times 4$  target matrix. Chen's algorithm is chosen here because it is simple and well recognized [13]. Normally, 8 4-point 1-D IDCTs are needed to accomplish the task (with very complicated data reordering, 4 4-point IDCTs are enough [8]). However, in our case at most 3 4-point IDCTs are necessary (1 along the rows (or columns) of the target matrix to get an intermediate matrix and 2 along the columns (or rows) of the intermediate matrix to get a  $2 \times 2$  submatrix of the partial contribution). The other three  $2 \times 2$  submatrices can be derived automatically by the use of mirror effect. Furthermore, each 4-point IDCT can be done efficiently since only one input data out of 4 is nonzero. Consider the signal flowgraph for a 4-point IDCT shown in Figure 2(a). The outputs of the 4-point IDCT (denoted by  $f_0$ ,  $f_1$ ,  $f_2$ , and  $f_3$ ) can be treated as linear combinations of the 4 inputs (denoted by  $F_0$ ,  $F_1$ ,  $F_2$ , and  $F_3$ ). Since only one of the inputs is nonzero, Figures 2(a) and 2(b) are functionally equivalent. The signal flowgraph in Figure 2(b) can be further simplified by retaining only two of the four outputs ( $f_0$  and  $f_1$ ) as shown in Figure 2(c) because the other two outputs can be derived by the use of mirror effect. Since the reducible property is separable, it can be used here to further reduce some of the subgraphs in Figure 2(c). Specifically, the subgraphs with input  $F_0$  and  $F_2$  are reducible. The



(a)



(b)

Figure 2. (a) A 4-point IDCT (b) signal decomposition of (a)  
(c) simplified version of (b) (d) irreducible subgraphs of (a)

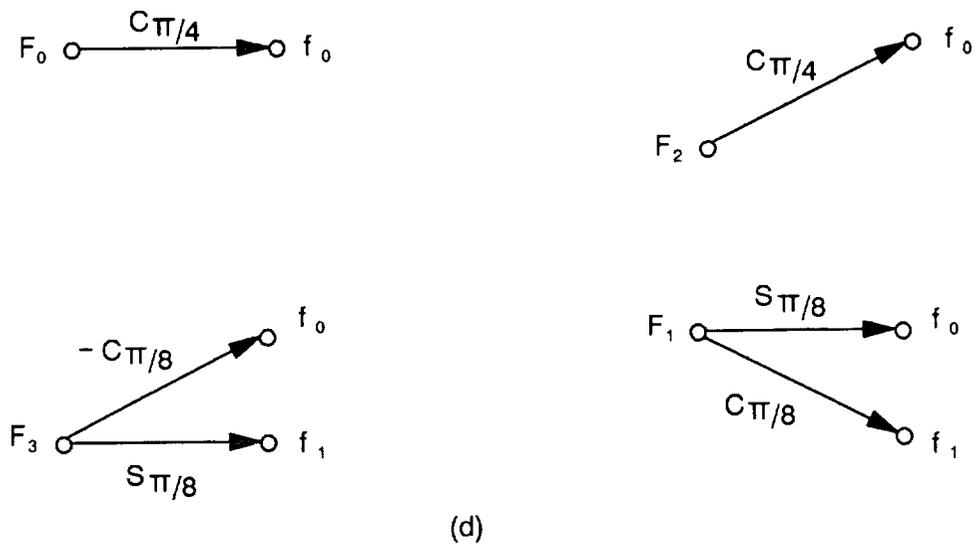
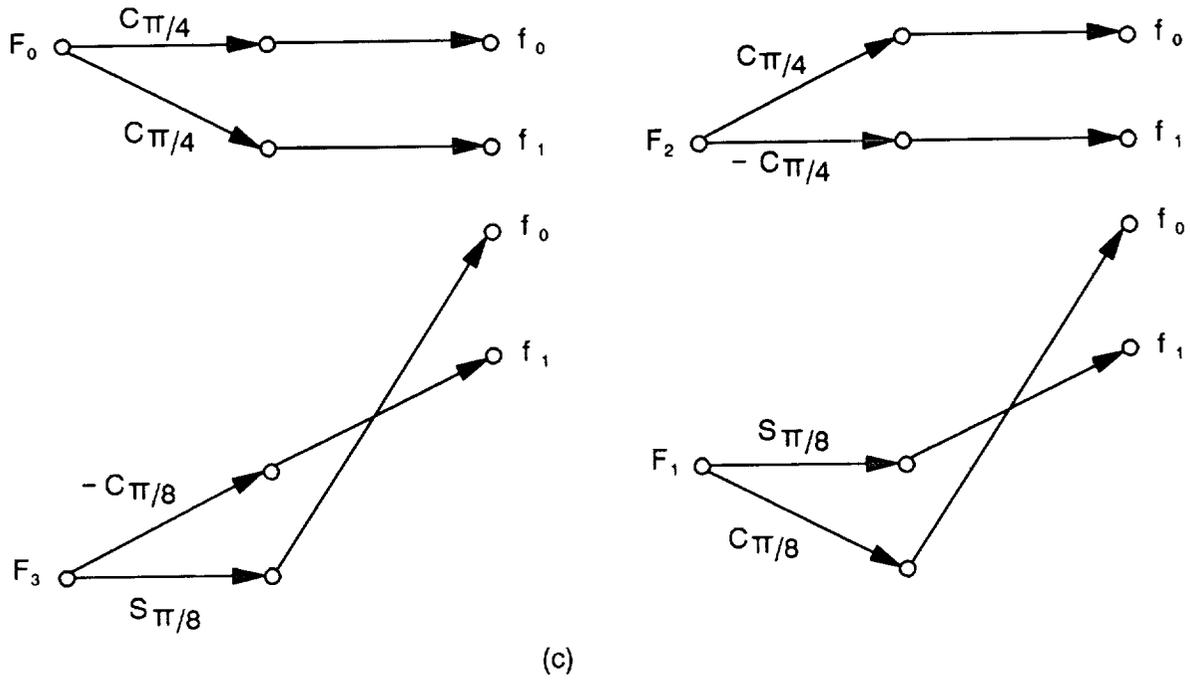


Figure 2. (Continued)

final irreducible subgraphs are shown in Figure 2(d). For convenience, the subgraphs shown in Figure 2(d) are said to be in their primitive forms. In other words, they can not be reduced or simplified any more. The above procedure can be extended easily to 8-point or higher order cases.

The primitive subgraph with input  $F_u$  and the one with input  $F_v$  can be cascaded as a signal flowgraph to compute part of the partial contribution due to  $F_{uv}$ . The connection rule is: at each output of the first subgraph, the second subgraph is cascaded. Which subgraph should be the first is immaterial as far as the result is concerned. However, the computational complexity may be different.

The complexity of computing part of the partial contribution due to  $F_{uv}$  can be examined by checking the primitive subgraphs with input  $F_u$  and  $F_v$ . The two primitive subgraphs are cascaded in the way described earlier. If the first subgraph takes  $P$  multiplications and the other requires  $Q$  multiplications, then the total number of multiplications required to obtain part of the partial contribution would be  $P+PQ$  multiplications. Alternatively,  $P$  and  $Q$  are also the number of output nodes for one subgraph and another, respectively. So  $P$  and  $Q$  can be obtained by counting the number of output nodes of the irreducible subgraphs. Since  $P+PQ = P(1+Q)$ , a fast way to tell the required number of multiplications is to take the product of  $P$  and  $Q+1$ .  $P+PQ$  multiplications will also be the complexity to compute the full size partial contribution since no addition operations are involved and the expansion of partial contribution to its full size adds no complexity. Suppose the two subgraphs are cascaded in reverse order, the complexity becomes  $Q+QP$ . But  $Q+QP \neq P+PQ$  if  $P \neq Q$ . Thus, the order of the subgraphs is relevant to the complexity. If  $P < Q$ ,  $P+PQ$  is always smaller than  $Q+QP$ . Therefore, the order selection should be such that the first one requires less complexity than the second one. The numbers of multiplications required for different combinations of  $u$  and  $v$  are shown in Table 2 and Table 3 for  $4 \times 4$  and  $8 \times 8$ , respectively. Note that for  $u=0$  or  $M/2$  and  $v=0$  or  $N/2$ , no multiplication is required (except a left shift operation by one bit).

Table 2. The Number of Multiplications Associated with  $F_{uv}$  ( $4 \times 4$ )

	$v$	0	1	2	3
$u$	0	0	3	0	3
	1	3	6	3	6
	2	0	3	0	3
	3	3	6	3	6

Table 3. The Number of Multiplications Associated with  $F_{uv}$  ( $8 \times 8$ )

	$v$	0	1	2	3	4	5	6	7
$u$	0	0	5	3	5	0	5	3	5
	1	5	20	10	20	5	20	10	20
	2	3	10	6	10	3	10	6	10
	3	5	20	10	20	5	20	10	20
	4	0	5	3	5	0	5	3	5
	5	5	20	10	20	5	20	10	20
	6	3	10	6	10	3	10	6	10
	7	5	20	10	20	5	20	10	20

According to Table 3, we can estimate the average number of multiplication required for a nonzero element in the 8 x 8 case. Assume that the chance of a nonzero element falling in any u-v pair is equally likely. Then the average will be the 1/64 of the sum of all the numbers shown in Table 3. The result is 9.5 multiplications per nonzero element. This is about 7 to 20 times faster than the naive approach mentioned earlier. Similarly, from Table 2, we will get 3 multiplications per nonzero element for the 4 x 4 case.

The update of a goal matrix is straightforward because it involves only additions of the corresponding elements in each partial contribution. The total number of additions is  $(n-1)MN$  for the update of the goal matrix, where n is the number of nonzero elements in an  $M \times N$  target matrix.

### 5. Performance Comparison of the Approaches

The advantages of the proposed approach are as follows:

- (1) It has essentially no delay time and computational redundancy.
- (2) It allows any scanning or transmission patterns, including the zig-zag scanning pattern. Note that the zig-zag scanning pattern is generally good for many images. However, a better or optimal scanning pattern for a particular image may deviate from the zig-zag scanning pattern [14]. Furthermore, it can be different from one image to another. Therefore, it is critical to be adaptive to different scanning patterns.
- (3) In inverse PIT, it has lower computational complexity than traditional fast algorithms.

Both advantages 1 and 2 are due to the separate processing of the element in the input matrix. The performance comparison of approach 1, approach 2, and the proposed approach are summarized in Table 4.

Table 4. Performance Comparison of the Approaches

	Approach 1	Approach 2	Proposed
Scanning Pattern	Adaptive	Fixed and Poor	Adaptive
Delay Time	High	Low	Lowest
Computational Redundancy	High	Low	Lowest
# of Multiplication*	384	114	95
# of Additions*	1864	740	576

\*The fast algorithm in [8] is used for approaches 1 and 2. The number of nonzero elements in an 8 x 8 target matrix is assumed to be 10 on the average. In addition, 4 stages of image reconstruction are assumed.

### 6. Conclusion

This paper presents a new and promising solution to the problem of heavy computation of IDCTs in inverse PIT process. When approach 2 was shown to have poor and discontinuous

image build up in 1990 [10], the research in fast progressive reconstruction for transform domain PIT schemes seems to be hopeless. With the proposed approach, both the fast progressive reconstruction and the pleasant image build up can be achieved simultaneously. This is an encouraging result for the research of transform domain fast progressive reconstruction.

### References

- [1] N. Dal Degan, P. Migliorati, S. Pozzi, V. Trecordi, "IMAGINE: Effective Retrieval from a Remote Image Database," Proceedings of IEEE Global Telecommunications Conference, Vol. 1, pp. 306-311, 1990.
- [2] Julius T. Tou, Robert W. Depree, Paul B. Lin, "Telebrowsing of Science Information Via a Minicomputer," Proceedings of NSF Seminar on Scientific and Technical Information Dissemination, Washing, D.C., 1976.
- [3] Julius T. Tou, Robert W. Depree, "Telebrowsing System and Its Applications," Proceedings of the European Computing Congress, 1978.
- [4] Kou-Hu Tzou, "Progressive Image Transmission: A Review and Comparison of Techniques," Optical Engineering, Vol. 26, No. 7, pp. 581-589, July 1987.
- [5] K. R. Rao, P. Yip, Discrete Cosine transform -- Algorithms, Advantages, and Applications, Academic Press, Inc., 1990.
- [6] Rafael C. Gonzalez, Richard E. Woods, Digital Image Processing, Addison-Wesley, Reading, MA, June 1992.
- [7] Gregory K. Wallace, "The JPEG Picture Compression Standard," IEEE Transactions on Consumer Electronics, Vol. 38, No. 1, pp. xviii-xxxiv, Feb. 1992.
- [8] Nam Ik Cho, Sang Uk Lee, "Fast Algorithm and Implementation of 2-D Discrete Cosine Transform," IEEE Transactions on Circuits and Systems, Vol. 38, No. 3, pp. 297-305, March 1991.
- [9] Kei Takikawa, "Fast Progressive Reconstruction of a Transformed Image," IEEE Transactions on Information Theory, Vol. 30, No. 1, pp. 111-117, January 1984.
- [10] M. Miran, K. R. Rao, "Fast Progressive Reconstruction of Images Using the DCT," in Signal Processing V: Theories and Applications, pp. 897-900, edited by L. Torres, E. Masgrau, and M. A. Lagunas, Elsevier Science Publishers B. V., 1990.
- [11] A. G. Tescher, R. V. Cox, "An Adaptive Transform Coding Algorithm," International Conference on Communications, pp. 47-20 through pp. 47-23, Philadelphia, PA, July 14-16, 1976.
- [12] I. S. Duff, A. M. Erisman, J. K. Reid, Direct Methods for Sparse Matrices, Clarendon Press, Oxford, 1986.
- [13] Wen-Hsiung Cheng, C. Harrison Smith, S. C. Fralic, "A Fast Computational Algorithm for the Discrete Cosine Transform," IEEE Transactions on Communications, Vol. 25, No. 9, pp. 1004-1009, September 1977.
- [14] Bowonkoon Chitprasert, K. R. Rao, "Human Visual Weighted Progressive Image Transmission," IEEE Transactions on Communications, Vol. 38, No. 7, pp. 1040-1044, July 1990.

## A SURVEY OF QUALITY MEASURES FOR GRAY SCALE IMAGE COMPRESSION

Ahmet M. Eskicioglu and Paul S. Fisher  
University of North Texas  
Department of Computer Science  
P.O. Box 13886  
Denton, TX 76203, USA  
E-mail: eskiciog@ponder.csci.unt.edu  
or fisher@gab.unt.edu

**Abstract.** Although a variety of techniques are available today for gray-scale image compression, a complete evaluation of these techniques cannot be made as there is no single reliable objective criterion for measuring the error in compressed images. The traditional subjective criteria are burdensome, and usually inaccurate or inconsistent. On the other hand, being the most common objective criterion, the mean square error (MSE) does not have a good correlation with the viewers' response. It is now understood that in order to have a reliable quality measure, a representative model of the complex human visual system is required. In this paper, we survey and give a classification of the criteria for the evaluation of monochrome image quality.

### 1. Introduction

There is an ever increasing demand for transmission and storage of vast amounts of information in data processing environments today. To reduce the large costs involved, data compression is a widely accepted tool which aims at minimizing the amount of data to be stored or transmitted. A variety of data compression techniques have been developed in the past few decades for different types of industrial, commercial, and educational applications. These techniques can be classified into two major categories: Lossless (exact) and lossy (inexact) [1, 2, 3]. Lossless compression is concerned with reconstructing an exact replica of the original input data stream. It is essentially used in text compression where no loss can be tolerated. Disastrous results may be encountered for even a single bit of loss in, for example, program files or database records. The techniques in this category typically reduce text size 40 to 80%, while those developed for specific applications may achieve compression over 90%. Lossy data compression causes some amount of loss which is considered to be a concession for a drastic increase in compression. Lossy compression techniques are effective and appropriate primarily for digitized voice and images for two reasons: Firstly, huge volumes of voice and images are normally generated in a typical application and, secondly, digital representation of analog signals is only an approximation, introducing a certain loss to begin with.

Numerous image compression techniques [2-6] exist today with the common goal of reducing the number of bits needed to store or to transmit images. The efficiency of a compression algorithm is generally measured using three criteria:

- 1) compression amount,
- 2) implementation complexity, and
- 3) resulting distortion.

The amount of compression can readily be obtained using several definitions, among which there are compression ratio, figure of merit, and compression percentage. Algorithmic complexity, on the other hand, can be measured by considering the data structures as well as the type and number

of operations required. The difficulty in evaluating a lossy compression algorithm comes from the fact that there is no reliable and consistent measure for determining the magnitude of distortion resulting from the loss. In other words, we lack a useful and practical measure for image quality assessment! Such a measure is not only needed for comparing images produced by different techniques, but it is also instrumental in designing image processing/compression algorithms.

In this paper, we survey the criteria available for the evaluation of monochrome image quality. In spite of the fact that some of the measures found in the literature have specifically been used for rating the performance of image processing systems, they are applicable in evaluating compression algorithms equally well.

## 2. Image Quality Measures

It is possible to classify image quality criteria as given in Figure 1.

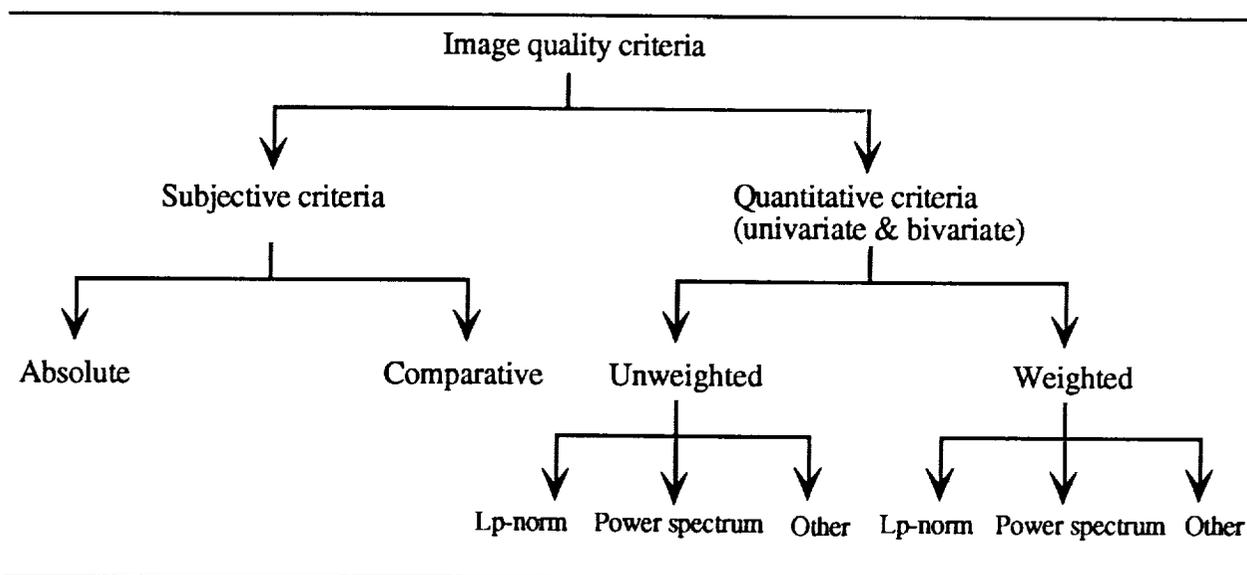


Figure 1. Classification of Image Quality Criteria

### 2.1 Subjective Criteria

As the final user of images are humans, the most reliable and commonly used assessment of image quality is the subjective rating by human observers. Both expert and nonexpert observers are used in experiments; nonexperts represent the average viewer while experts are believed to be able to give better, more 'refined' assessments of image quality since they have been trained and are familiar with images and their distortions.

In absolute evaluation, the observers view an image and assess its quality by assigning to it a category in a given rating scale, whereas in comparative evaluation, a set of images are ranked from best to worst by the observers. The rating scales that appear in the relevant literature [5, 12, 14, 15, 19] are listed in Table 1.

**Table 1. Rating Scales Used in Subjective Evaluation**

<p>A.</p> <ol style="list-style-type: none"> <li>5. Excellent</li> <li>4. Good</li> <li>3. Fair</li> <li>2. Poor</li> <li>1. Unsatisfactory (bad)</li> </ol>	<p>B.</p> <ol style="list-style-type: none"> <li>7. Best</li> <li>6. Well above average</li> <li>5. Slightly above average</li> <li>4. Average</li> <li>3. Slightly below average</li> <li>2. Well below average</li> <li>1. Worst</li> </ol>	<p>C.</p> <ol style="list-style-type: none"> <li>1. Not noticeable (perceptible)</li> <li>2. Just noticeable (perceptible)</li> <li>3. Definitely noticeable (perceptible) but only slight impairment</li> <li>4. Impairment not objectionable</li> <li>5. Somewhat objectionable</li> <li>6. Definitely objectionable</li> <li>7. Extremely objectionable</li> </ol>										
<p>D.</p> <ol style="list-style-type: none"> <li>3 Much better</li> <li>2 Better</li> <li>1 Slightly better</li> <li>0 Same</li> <li>-1 Slightly worse</li> <li>-2 Worse</li> <li>-3 Much worse</li> </ol>	<p>E.</p> <ol style="list-style-type: none"> <li>5. Imperceptible</li> <li>4. Perceptible but annoying</li> <li>3. Slightly annoying</li> <li>2. Annoying</li> <li>1. Very annoying</li> </ol>	<p>F.</p> <table style="border: none;"> <tr> <td>10, 9</td> <td>Very good</td> </tr> <tr> <td>8, 7</td> <td>Good</td> </tr> <tr> <td>6, 5, 4</td> <td>Fair</td> </tr> <tr> <td>3, 2</td> <td>Bad</td> </tr> <tr> <td>1, 0</td> <td>Very bad</td> </tr> </table>	10, 9	Very good	8, 7	Good	6, 5, 4	Fair	3, 2	Bad	1, 0	Very bad
10, 9	Very good											
8, 7	Good											
6, 5, 4	Fair											
3, 2	Bad											
1, 0	Very bad											

The mean rating of a group of observers who join the evaluation is usually computed by

$$R = \left( \sum_{k=1}^n s_k n_k \right) / \left( \sum_{k=1}^n n_k \right),$$

where  $s_k$  = the score corresponding to the  $k$ th rating,  $n_k$  = the number of observers with this rating, and  $n$  = the number of grades in the scale.

Bubble sort [5, 11, 22] is another technique used in image rating. With this technique, the subject compares two images A and B from a group and determines their order. Assuming that the order is AB, he/she takes a third image and compares it with B to establish the order ABC or ACB. If the order is ACB, then another comparison is made to determine the new order. The procedure continues until all the images have been used, allowing the best pictures to bubble to the top if no ties are accepted.

It is important to note that the results of subjective rating are affected by a number of factors including

- a) type and range of images,
- b) level of expertise of the observers, and
- c) experimental conditions.

If standards can be established for these factors, the results obtained in different locations and at different times may then become comparable.

## 2.2 Quantitative Criteria

Quantitative measures for image quality can be divided into two classes: Univariate and bivariate [19]. A univariate measure assigns to a single image a numerical value based upon measurements of the image field, and a bivariate measure is a numerical comparison between two images.

Fidelity measurements are usually made using an array of discrete image samples, although a continuous image field can also be generated by two-dimensional interpolation of the sample array if the overhead is justified. Image error measures can be defined in either spatial or frequency domain.

Denoting the samples on the original image field as  $F(j,k)$ , a spatial domain, univariate quality rating may be expressed in general as

$$Q = \sum_{j=1}^M \sum_{k=1}^N O\{F(j,k)\}$$

for  $N \times M$  samples, where  $O\{\cdot\}$  is some operator.

Bivariate measures are more frequently used in image quality measurement. If  $\hat{F}(j,k)$  denotes the samples on the degraded image field, a number of measures can be established to determine the closeness of the two image fields. The alternatives are listed below [5, 9, 12, 19, 22-25].

$$(i) \quad L_p = \left\{ (1/MN) \sum_{j=1}^M \sum_{k=1}^N |F(j,k) - \hat{F}(j,k)|^p \right\}^{1/p}$$

A major class of bivariate error measures is based on the  $L_p$ -norm. The factor  $p$  determines the relative significance of errors of different magnitudes.  $L_1$  is the average absolute error and  $L_2$  is the commonly used root mean square error (RMSE). As the value of  $p$  is increased, a greater relative emphasis is given to large errors in the image.

(ii) Low order moment of a power spectrum.

$$(iii) \quad K = \sum_{j=1}^M \sum_{k=1}^N F(j,k) \hat{F}(j,k)$$

This measure is obtained by discretizing the continuous cross-correlation function. It may be normalized by the reference image energy to give unity as the peak correlation:

$$NK = \frac{\sum_{j=1}^M \sum_{k=1}^N F(j,k) \hat{F}(j,k)}{\sum_{j=1}^M \sum_{k=1}^N [F(j,k)]^2}$$

(iv) Correlation quality:

$$CQ = \frac{\sum_{j=1}^M \sum_{k=1}^N F(j,k) \hat{F}(j,k)}{\sum_{j=1}^M \sum_{k=1}^N F(j,k)}$$

(v) Structural content:

$$SC = \frac{\sum_{j=1}^M \sum_{k=1}^N [F(j,k)]^2}{\sum_{j=1}^M \sum_{k=1}^N [\hat{F}(j,k)]^2}$$

(vi) Normalized absolute error between the reference and degraded image fields:

$$NAE = \frac{\sum_{j=1}^M \sum_{k=1}^N |O\{F(j,k)\} - O\{\hat{F}(j,k)\}|}{\sum_{j=1}^M \sum_{k=1}^N |O\{F(j,k)\}|}$$

(vii) Normalized mean square error:

$$NMSE = \frac{\sum_{j=1}^M \sum_{k=1}^N [O\{F(j,k)\} - O\{\hat{F}(j,k)\}]^2}{\sum_{j=1}^M \sum_{k=1}^N [O\{F(j,k)\}]^2}$$

(viii) Peak mean square error:

$$PMSE = \frac{(1/MN) \sum_{j=1}^M \sum_{k=1}^N [O\{F(j,k)\} - O\{\hat{F}(j,k)\}]^2}{A^2}$$

where A represents the maximum value of O{F(j,k)}.

The definitions used for the operator O{·} in (vii) and (viii) are

- (a) F(j,k)
- (b) [F(j,k)]<sup>v</sup> (Power law)

- (c)  $k_1 \log_b [k_2 + k_3 F(j,k)]$  (Logarithmic)
- (d)  $[F(x,y) \otimes H(x,y)] \delta(x-j\Delta x, y-k\Delta y)$  (Convolution)

(ix) Laplacian mean square error:

$$\text{LMSE} = \frac{\sum_{j=1}^{M-1} \sum_{k=2}^{N-1} [O\{F(j,k)\} - O\{\hat{F}(j,k)\}]^2}{\sum_{j=1}^{M-1} \sum_{k=2}^{N-1} [O\{F(j,k)\}]^2}$$

where  $O\{F(j,k)\} = F(j+1, k) + F(j-1, k) + F(j, k+1) + F(j, k-1) - 4F(j,k)$

In many applications, the mean square error (however it is defined) is often expressed in terms of a signal-to-noise ratio defined in decibels.

(x) Image fidelity:

$$\text{IF} = 1 - \frac{\sum_{j=1}^M \sum_{k=1}^N [F(j,k) - \hat{F}(j,k)]^2}{\sum_{j=1}^M \sum_{k=1}^N [F(j,k)]^2}$$

(xi) Difference  $[j,k] = F(j,k) - \hat{F}(j,k)$

(xii)  $\sum_{j=1}^M \sum_{k=1}^N \text{Difference } [j,k] / MN$

(xiii)  $\text{Max}\{|\text{Difference}[j,k]|\}$

(xiv) Histogram of the compression error (constructed by plotting the number of x's versus x for all values of x found in the difference matrix).

(xv) Hosaka plots

(xvi) Sensitivity and predictive value positive curves

(xvii) Rate-distortion curves.

It is reported that image quality assessment can be improved by incorporating into the evaluation process some model of the HVS. The HVS is incorporated into the quality measure using two distinct approaches. In the first approach, the  $L_p$  norm (or one of its variants) is employed attaching a weight to the image samples either in the spatial or frequency domain. The second approach is concerned with weighting the digital image power spectrum.

In one of the earliest studies, the transformation

$$O\{\cdot\} = H_L(x,y) \otimes O_N\{\cdot\}$$

is used on both the continuous image field  $F(x,y)$  and the degraded image field  $\hat{F}(x,y)$  before applying the integral square error, where the impulse response  $H_L(x,y)$  represents the lateral inhibition process, and the point nonlinearity  $O_N\{\cdot\}$  models the response of the eye's photoreceptors [11]. In the Fourier domain  $H_L$  is defined as

$$a \left[ c + \left( \frac{\omega}{\omega_0} \right)^{k_1} \right] \exp \left\{ - \left( \frac{\omega}{\omega_0} \right)^{k_2} \right\},$$

where  $\omega = (\omega_1 + \omega_2)^{1/2}$ , and  $O\{\cdot\} = \{\cdot\}^{1/3}$  is chosen. The experiments show that  $a = 2.6$ ,  $c = 0.0192$ ,  $\omega_0 = 1/0.114$ ,  $k_1 = 1$  and  $k_2 = 1.1$  are the suitable parameter values.

In another study [12] to find an objective measure which closely mirrors the performance of the human viewer, the error measure

$$E_p = \left\{ \frac{1}{m} \sum_{i=1}^m |e_i|^p \right\}^{1/p}$$

where  $m$  = number of picture elements (pels) in a picture,  $e_i = x_i - \hat{x}_i$ ,  $x_i$  = the value of the pel in the original picture and  $\hat{x}_i$  = the value of the pel in the distorted picture, is tried for  $p = 1, 2, 3, 4, 6$ . The conclusion is that  $E_p$  is a very good estimate of impairment rating where the type of distortion is additive white noise. In the same study, another measure of picture impairment is obtained using

$$EM_p = \left\{ \frac{1}{m} \sum_{i=1}^m |e_i|^p / W_i \right\}^{1/p}$$

to reflect the masking effect of the signal.  $W_i$  denotes the value of the weighting function at pel  $i$  and is derived from an activity function that is a measure of the variability of the signal in the neighborhood of pel  $i$ . Three different forms of activity functions are studied:

$A_{max}$ : measures the maximum signal change between any pair of pels in a neighborhood consisting of the pel being evaluated plus the eight surrounding pels.

$A_{av}$ : sums the deviations of the same neighborhood of points from the neighborhood average  $\bar{x}$

$A_{df}$ : provides the weighted sum of the magnitude of the surrounding element difference (slope) in both the horizontal and vertical directions.

In all three cases  $W_i$  is obtained from  $A_i$  so as to span a range from 1.0 to 10.0. There is also an attempt in [12] to obtain a local measure of image quality. Relying on the postulate that the viewer rates the image by some weighted average of the worst two or three patches, Limb divides the image into a rectangular array of squares and calculates a local measure for each square with and without masking. He also tries the formula

$$EM_p = \left\{ \frac{1}{m} \sum_{i=1}^m |e_i/W_i|^p \right\}^{1/p}$$

in his local error analysis. The quantitative model that Limb uses for the human viewer includes some error filtering as well. Comparison of the simple RMSE as a measure of image quality with the best error measure predictions of the model shows that RMSE performs surprisingly well. This results, Limb explains, from the fact that in most distorted images, quality is determined mainly by the visibility of distortion in flat areas where it is more visible and consequently the effects of masking have little effect. For images where distortion is greater at edges, however, the RMSE is claimed to be less satisfactory.

The results of a subjective evaluation on twelve versions of a black and white image and the rank ordering obtained with three computational measures are presented by Hall [22]. He compares the performance of the measures NMSE, LMSE, and PMSE, which are defined for an  $N \times N$  discrete image as

$$NMSE = \frac{\sum_{m=1}^N \sum_{n=1}^N [f(m,n) - \hat{f}(m,n)]^2}{\sum_{m=1}^N \sum_{n=1}^N [f(m,n)]^2}$$

$$LMSE = \frac{\sum_{m=2}^{N-1} \sum_{n=2}^{N-1} [G(m,n) - \hat{G}(m,n)]^2}{\sum_{m=2}^{N-1} \sum_{n=2}^{N-1} [G(m,n)]^2}$$

where  $G(m,n) = f(m+1,n) + f(n-1,n) + f(m,n+1) + f(m,n-1) - 4f(m,n)$

$$PMSE = \frac{\sum_{m=1}^N \sum_{n=1}^N [z(m,n) - \hat{z}(m,n)]^2}{\sum_{m=1}^N \sum_{n=1}^N [z(m,n)]^2},$$

where  $z(m,n)$  and  $\hat{z}(m,n)$  are given by

$$z(m,n) = \ln[f(m,n)] \otimes h_{bp}(m,n)$$

and

$$\hat{z}(m,n) = \ln[\hat{f}(m,n)] \otimes h_{bp}(m,n)$$

The function  $h_{bp}(m,n)$  is a rectangular coordinate form of the point spread function of the HVS. In his comparison, Hall finds that the correlation between PMSE and the subjective ranking (obtained by using bubble sort) of the data set is higher than that of NMSE and LMSE.

Neill [8] arrives at a quality measure in the 2-D discrete Fourier spatial frequency domain. This measure is expressed as

$$K^{-1} \sum_{i=1}^B W_i \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} H^2(r) [F_i(u,v) - \hat{F}_i(u,v)]^2,$$

where  $B$  = number of subimage blocks in scene,

$K$  = normalization factor such as total energy,

$H(r)$  = rotationally symmetric spatial frequency response of HVS,  $r = \sqrt{u^2 + v^2}$ ,

$F_i, \hat{F}_i$  = Fourier transform of unprocessed and processed subimage  $i$ , respectively,

$M, N$  = number of Fourier coefficients + 1, in orthogonal  $u, v$  directions,

$W_i$  = subimage  $i$  structure weighting factor, proportional to subimage's intensity level variance.

Using  $H(r) = (0.2 + 0.45r)e^{-0.18r}$ , he then constructs the function

$$|A(r)| H(r) = \begin{cases} 0.05r^{0.554}, & \text{for } r < 7 \\ e^{-9 \left[ \log_{10} r - \log_{10} 9 \right]^2 2.3} & \text{for } r \geq 7 \end{cases}$$

for dealing with image cosine transforms instead of image Fourier transforms. Finally, he argues that (i) combining the HVS model with the image cosine transform will result in better performance in image compression and image quality assessment applications, and (ii) performance in quality assessment should also be enhanced by inclusion of the subimage structure weighting.

Marmolin [9] addresses the question of using the mean squared error (MSE) measure as a quality criterion in image processing, and evaluates the predictive power of

$$E = \left[ \frac{1}{n} \sum_{i=1}^n |D_i|^p \right]^{1/p},$$

$$D_i = a_i - g(x_i - y_i)$$

where  $g$  = some processing function that determines the visibility of the error,  $a_i$  = a weight related to the informative value of pixel  $i$ , and  $p$  = a factor that determines the relative importance of small and large errors,  $x_i$  = the gray level of pixel  $i$  in the original image,  $y_i$  = the gray level of pixel  $i$  in the processed image. He investigates the performance of different definitions for  $D_i$ , and compares them to that of the mean squared error

$$\text{MSE} = \left[ \frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2 \right]^{1/2}$$

The results obtained indicate that MSE is an unsatisfactory measure of perceived similarity, and that no measure is valid for each image set used.

Saghri, Cheatham, and Habibi [10] state that once an image  $U(x,y)$  and its reproduction have been subjected to the HVS model, then the mean square error

$$d(U,U') = \frac{1}{N} \iint [U(x,y) - U'(x,y)]^2 dx dy,$$

where  $N$  is the image area or the number of pixels, may be considered as a meaningful measure of image quality. Adopting the approach of Mannos and Sakrison, they use in their HVS model

$$f(u) = u^{0.33}$$

where  $u$  is the pixel intensity, and

$$A(f_r) = \left[ 0.2 + 0.81 \left( \frac{f_r}{5.55} \right) \right] \exp \left[ - \left( \frac{f_r}{5.55} \right) \right],$$

where  $f_r = (f_x^2 + f_y^2)^{1/2}$ . The corrections (developed by Nill)

$$C(f_r) = \left( \frac{1}{4} + \frac{1}{2} \left\{ \log_e \left[ \frac{2\Pi f_r}{\alpha} + \left( \frac{4\Pi^2 f_r^2}{\alpha} + 1 \right)^{1/2} \right] \right\}^2 \right)^{-1}$$

to the HVS model of  $A(f_r)$  is then added to give the DCT version

$$A_{DCT} = A(f_r)C(f_r).$$

As an alternative to the MSE, the authors propose the so-called information content (IC). The IC of an image for a given resolution is defined as the sum of the magnitudes of its DCT spectral components after they have been appropriately normalized based on HVS sensitivity models for that particular resolution. The plot of IC versus the resolution provides some insight into the quality of a given image. The preliminary results are reportedly promising, but much more experimentation is needed to adjust the numerous parameters of the system for highest achievable correlation with the subjective measure.

The work by Ngan, Leong, and Singh [16] describes an adaptive cosine transform coding scheme for color images. The cosine transform coefficients are weighted by the HVS function given by Nill to generate the coefficients in perceptual domain. To determine the parameters of the HVS filter

$$H(\omega) = (a+b\omega) \exp(-c\omega)$$

plots of SNR versus peak frequency are used. The SNR is defined by

$$SNR = -10 \log_{10} \left[ \frac{1}{(512)^2} \sum_{j=0}^{511} \sum_{k=0}^{511} \frac{[f(j,k) - \hat{f}(j,k)]^2}{(255)^2} \right],$$

where  $f(j,k)$  and  $\hat{f}(j,k)$  are the original and reconstructed pixels, respectively. Their results show that the subjective quality of the reconstructed images at a bit rate of 0.4 bit/pixel or a compression ratio of 60:1 is very good.

Khafizov, Fisher, and Kiselyov [18] propose a new approach to simulate human visual perception in order to devise a tool for measuring distance between images. Defining the error matrix by

$$E = X - Y,$$

where  $X$  and  $Y$  are the two images to be compared, they renormalize each error in  $E$  with respect to other errors. Renormalization is the core of their method and it produces a new re-estimated error matrix  $E'$ . Once  $E'$  is obtained, they compute the  $L_1$ -norm of  $E'$  as the distance between  $X$  and  $Y$ . In the case when there are only two errors  $e$  and  $z$  in  $E$ , the formula

$$e'(z) = \frac{3+a^s}{z(1+a^s)} (e+\tilde{z}), \text{ where } \tilde{z} = \begin{cases} z, & ez > 0 \\ 2e-z, & ez < 0 \end{cases}$$

where  $a$  = some positive constant and  $s$  = distance between  $e$  and  $z$ , is used for re-estimating the error  $e$  with respect to error  $z$ . The generalization to an arbitrary case is immediate. The experiments presented demonstrate the inconsistency of the conventional RMSE together with the success in simulating visual human perception.

Nil and Bouzas [17] present an objective, quantitative image quality measure based on the digital image power spectrum of normally acquired arbitrary scenes. Using polar coordinates  $\rho, \theta$  the image quality measure is derived from the normalized 2-D power spectrum  $P(\rho, \theta)$  weighted by the square of the modulation transfer function of the human visual system  $A^2(T\rho)$ , the directional scale of the input image  $S(\theta_1)$ , and the modified Wiener noise filter  $W(\rho)$ :

$$IQM = \frac{1}{M^2} \sum_{\theta=-180}^{180} \sum_{\rho=0.01}^{0.5} S(\theta_1) W(\rho) A^2(T\rho) P(\rho, \theta),$$

where  $M^2$  = number of pixels. In its application, a previously constructed modulation transfer function [8] is used for the HVS. The authors point out that the power spectrum approach does not require use of designed quality assessment targets or reimagining the same scene for comparison purposes. Experimental verification indicates good correlation of this objective quality measure with visual quality assessments.

### 3. Conclusions

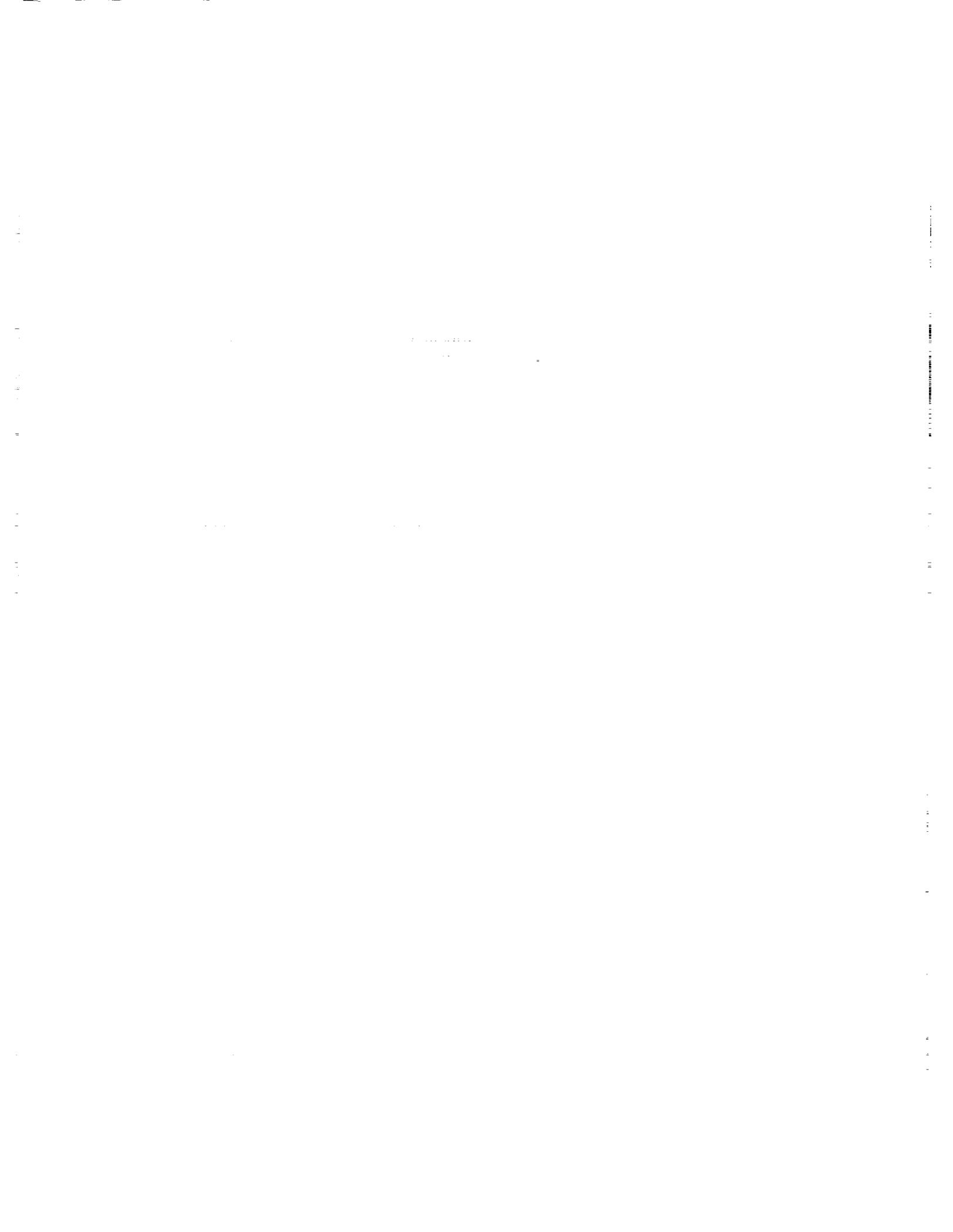
Traditionally, the most reliable way of measuring image quality has been the subjective evaluation by human observers. Because of the inherent difficulties associated with this approach, much attention has been focused on the development of quantitative techniques for quick and objective measurement. The image quality measure that has been commonly used in digital image compression is the mean square error (MSE) between the original image and the reconstructed image. It is now a well-known fact, however, that the MSE and its variants do not correlate reasonably well with subjective quality measures [4, 5, 7-10, 21]. A major portion of recent research is, therefore, directed towards incorporating human visual system (HVS) models into image quality measures. This is not a trivial task because the human visual system is too complex and an accurate model cannot presently be developed. Nevertheless, a number of experiments with simplified models indicates that the inclusion of a model for the HVS generally produces results that are in better correlation with the perceived image quality [4, 7, 8, 10-18, 22]. The trial models take into consideration various recognized characteristics of the HVS, and usually have both linear

and nonlinear parts. As we have a better understanding of the psychophysical phenomena concerning the human vision, we will be able to develop more accurate models which, in turn, will lead to results closer to the human response.

### References

- [1] T.C. Bell, J.G. Cleary, and I.H. Witten, *Text Compression*, Prentice-Hall, Inc., USA, 1990.
- [2] J.A. Storer (Editor), *Image and Text Compression*, Kluwer Academic Publishers, USA, 1992.
- [3] M. Nelson, *The Data Compression Book*, M&T Publishing, Inc., 1992.
- [4] A.K. Jain, "Image Data Compression: A Review," *Proceedings of the IEEE*, Vol. 69, No. 3, pp. 349-389, March 1981.
- [5] A.K. Jain, *Fundamentals of Digital Image Processing*, Prentice-Hall, Inc., USA, 1989.
- [6] M.P. Ekstrom (Editor), *Digital Image Processing Techniques*, Academic Press, Inc., USA, 1984.
- [7] D.J. Granrath, "The Role of the Human Visual Models in Image Processing," *Proceedings of the IEEE*, Vol. 69, No. 5, pp. 552-561, May 1981.
- [8] N.B. Nill, "A Visual Model Weighted Cosine Transform for Image Compression and Quality Assessment," *IEEE Transactions on Communications*, Vol. COM-33, No. 6, pp. 551-557, June 1985.
- [9] H. Marmolin, "Subjective MSE Measures," *IEEE Transactions on Systems, Man, and Cybernetics*, Vol. SMC-16, No. 3, pp. 486-489, May/June 1986.
- [10] J.A. Saghri, P.S. Cheatham, and A. Habibi, "Image Quality Measure Based on a Human Visual System Model," *Optical Engineering*, Vol. 28, No. 7, pp. 813-818, July 1989.
- [11] J.L. Mannas, "The Effects of a Visual Fidelity Criterion on the Encoding of Images," *IEEE Transactions on Information Theory*, Vol. IT-20, No. 4, pp. 525-536, July 1974.
- [12] J.O. Limb, "Distortion Criteria of the Human Viewer," *IEEE Transactions on Systems, Man, and Cybernetics*, Vol. SMC-9, No. 12, pp. 778-793, December 1979.
- [13] D.J. Sakrison, "On the Role of the Observer and a Distortion Measure in Image Transmission," *IEEE Transactions on Communications*, Vol. COM-25, No. 11, pp. 1251-1267, November 1977.
- [14] A.N. Netravali and J.O. Limb, "Picture Coding: A Review," *Proceedings of the IEEE*, Vol. 68, No. 3, pp. 366-406, March 1980.
- [15] F.X.J. Lukas and Z.L. Budrikis, "Picture Quality Prediction Based on a Visual Model," *IEEE Transactions on Communications*, Vol. COM-30, No. 7, pp. 1679-1692, July 1982.

- [16] K.N. Ngan, K.S. Leong, and H. Singh, "Adaptive Cosine Transform Coding of Images in Perceptual Domain," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, Vol 37, No. 11, pp. 1743-1750, December 1989.
- [17] N.B. Nill and B.H. Bouzas, "Objective Image Quality Measures Derived From Digital Image Power Spectra," *Optical Engineering*, Vol. 31, No. 4, pp. 813-825, April 1992.
- [18] F.T. Khafizor, P.S. Fisher, and O. Kiselyov, "A Note on Comparing Images," submitted to *IEEE Computing*, 1992.
- [19] W.K. Pratt, *Digital Image Processing*, John Wiley and Sons, Inc., USA, 1978.
- [20] N.B. Nill, "Scene Power Spectra: The Moment as an Image Quality Merit Factor," *Applied Optics*, Vol. 15, No. 11, pp. 2846-2854, November 1976.
- [21] D.R. Ahlgren, J. Crosbie, and D. Erigat, "Compression of Digitized Images for Transmission and Storage Applications," *Proceedings of SPIE*, Vol. 901, pp. 105-113, 1988.
- [22] C.F. Hall, "Subjective Evaluation of a Perceptual Quality Metric," *Proceedings of SPIE*, Vol. 310, pp. 200-204, 1981.
- [23] H.L. Snyder, "Image Quality: Measures and Visual Performance," *Flat-Panel Displays and CRTs*, L.E. Tannas, Jr., Ed., Van Nostrand Reinhold, New York, , pp. 70-90, 1985.
- [24] P.M. Farrelle, *Recursive Block Coding for Image Data Compression*, Springer-Verlag New York Inc., USA, 1990.
- [25] P.C. Cosman, C. Tseng, R.M. Gray, R.A. Olshen, L.E. Moses, H.C. Davidson, C.J. Bergin, and E.A. Riskin, "Tree-structured Vector Quantization of CT Chest Scans: Image Quality and Diagnostic Accuracy," Technical Report No. 157, Division of Biostatistics, Stanford University, Stanford, California, August 1992.



## Digital Mammography, Cancer Screening: Factors Important for Image Compression

Laurence P. Clarke<sup>1</sup>, G. James Blaine<sup>2</sup>, Kunio Doi<sup>3</sup>, Martin J. Yaffe<sup>4</sup>, Faina Shtern<sup>5</sup>, and G. Stephen Brown<sup>5</sup>; *NCI/NASA Working Group on Digital Mammography*  
 Daniel L. Winfield; *Research Triangle Institute, Research Triangle Park, NC 27709-2194*  
 Maria Kallergi<sup>6</sup>; *Department of Radiology, University of South Florida, Tampa, FL 33612-4799*

**Abstract.** The use of digital mammography for breast cancer screening poses several novel problems such as development of digital sensors, computer assisted diagnosis (CAD) methods for image noise suppression, enhancement, and pattern recognition, compression algorithms for image storage, transmission, and remote diagnosis. X-ray digital mammography using novel direct digital detection schemes or film digitizers results in large data sets and, therefore, image compression methods will play a significant role in the image processing and analysis by CAD techniques. In view of the extensive compression required, the relative merit of "visually lossless" versus lossy methods should be determined. A brief overview is presented here of the developments of digital sensors, CAD, and compression methods currently proposed and tested for mammography. The objective of the NCI/NASA Working Group on Digital Mammography is to stimulate the interest of the image processing and compression scientific community for this medical application and identify possible dual use technologies within the NASA centers.

### 1. Introduction

Mammography is widely accepted today as the most effective method of screening women for breast cancer [1,2,3]. Recent studies indicate that approximately 14,000 women's lives were saved in 1992 through mammographic screening programs. Digital mammography promises significant improvements over currently used screen-film mammography. Its aims are to integrate (a) solid state sensors for digital localized or full-view breast imaging to improve image quality and acquisition process; (b) computer algorithms for image enhancement and extraction of features such as calcifications and masses to assist the radiologist and improve screening and diagnosis; (c) image transmission and storage techniques for telemammography to improve patient care and efficient use of professional expertise [4,5]. The technical challenges in digital mammography are similar to those confronted in other scientific areas, e.g. space-based sensors and space-generated information, indicating significant possibilities for cross fertilization of ideas and dual use technologies. The NCI/NASA Working Group on Digital Mammography has been

---

<sup>1</sup> Center for Engineering and Medical Image Analysis (CEMIA), University of South Florida, Tampa, FL, 33612-4799.

<sup>2</sup> Mallinckrodt Institute of Radiology, Washington University Medical Center, St. Louis, MO 63110.

<sup>3</sup> Kurt Rossmann Laboratories for Radiologic Image Research, University of Chicago, Chicago, IL 60637.

<sup>4</sup> Medical Physics Research Group, Sunnybrook Health Science Centre and Departments of Medical Biophysics and Radiology, University of Toronto, North York, Ontario, Canada.

<sup>5</sup> Diagnostic Imaging Research Branch, National Cancer Institute, Bethesda, MD 20892.

<sup>6</sup> Consultant to Research Triangle Institute

established to explore application opportunities of such dual use technologies from aerospace to medicine [4]. In the following sections, the general requirements are presented in greater detail followed by a summary of the technologies identified until now by this Working Group within NASA Centers and Federal Laboratories as relevant to digital mammography and with a potential of having an impact to the problem.

## 2. Digital Mammography Systems

Conventional x-ray film-screen mammography is currently an accepted imaging modality for breast cancer detection. Nevertheless there are several technical factors limiting its performance [6]. These are the tradeoffs between contrast and exposure range inherent in the film-based system, the influence of film grain on image noise and the inefficiency of conventional methods for rejecting scattered x-rays. A direct method of mammographic image acquisition in digital form can overcome these limitations and provide improved detection of breast cancer. *The target specifications of a digital mammography detector are: (1) efficient absorption of the incident radiation beam; (2) linear response over a wide range of incident radiation intensity; (3) low intrinsic noise; (4) high spatial resolution; 50  $\mu\text{m}$  maximum; (5) high dynamic range; 4,000:1 minimum; (6) should accommodate at least an 18 x 24 cm field size; (7) allow an acceptable imaging time (1-7 s) and heat loading of the x-ray tube; (8) display capabilities: pixel matrices of 2k x 2k for soft copy and 4k x 4k for hard copy.*

Various configurations for image acquisition have been considered such as area, point, line and multiline systems. Each approach involves compromises between factors such as spatial resolution, imaging time, readout time, detector dynamic range and sensitivity, cost, susceptibility to artifacts, efficiency of scatter reduction, and available detector size.

Photostimulable phosphors with laser readout is an approach which has been successfully developed for general radiography and it is possible to extract the information from such systems in digital form. Currently, however, this technology does not provide adequate spatial resolution for mammography [7] and, because of inefficiencies in signal collection, may suffer from excessive image noise since the detector may not be x-ray quantum limited at mammographic energies. Interesting developments in area detectors are currently underway. They include large area charge-coupled devices (CCDs) [8], silicon or amorphous silicon [9,10], amorphous selenium [11], and improvements on photostimulable phosphors. One or more of these could play a future role in digital mammography.

An approach particularly attractive for digital mammography is through *area detectors*. Such an approach is convenient and makes efficient use of the heat loading applied to the x-ray tube. Unfortunately, area detectors which combine adequate spatial resolution and field coverage with good signal-to-noise characteristics do not currently exist. For example, simple coupling of a large-area phosphor to a small-area photodetector via lenses requires a large minification factor ( $M$ ). Because the efficiency of light transfer is approximately proportional to  $M^{-2}$ , this is an inefficient means of imaging which causes the system not to be x-ray quantum limited.

Although areas detectors probably present the most acceptable long range solution to digital mammography, it is not clear how long it will take to overcome the technical challenges to produce a practical clinical system of this type. It is, however, currently possible to meet the specifications described above using a scanned-beam method of image acquisition [12]. The superior efficiency of scatter reduction inherent in a scanning system compared to an area detector can provide advantages in terms of image SNR/radiation dose. For scanning systems,

x-ray tube heat loading is always a concern. Scanned point and single line systems are impractical for this reason.

At the University of Toronto (UT), a slot-beam imaging system for digital mammography, shown schematically in Figure 1, is currently under development [13]. The radiation beam forms a "slot" of dimensions approximately 24 cm by 4 mm. After transmission through the breast, x-rays are incident on a fluorescent phosphor, and the emitted light enters a fiber optic assembly consisting of two fan-shaped tapers. The end of each of the 2X demagnifying tapers is ground to a 45° angle at the detector input surface where the two tapers are fused together in a smooth joint without a line that is parallel to the scan direction. The output surface of each taper is mated to a CCD array. This arrangement provides a pixel size of 50 µm referred to the midplane of the breast.

The image is acquired using time delay integration (TDI), by scanning the fan x-ray beam and the slot detector across the breast in a direction parallel to the short dimension of the detector [14]. The 45° joint of the input surface of the fiber optic tapers and the TDI acquisition solves one of the major problems associated with modular detectors, which is the presence of artifacts at their junctions. The TDI motion will average the variations in signal along detector columns due to detector structure, including that due to the joint between modules, thereby avoiding disturbing artifacts. A secondary advantage of the TDI acquisition is that both dark current and detector uniformity corrections can be made by acquiring "image" data first without x-rays and then with a uniform x-ray exposure to the detector and sorting one offset correction and one gain factor for each of the 4096 detector columns.

Current state-of-the-art does not yet provide adequate "soft-copy" display resolution ( $\geq 4096 \times 4096$  pixels). A high resolution CRT display (2k x 2k pixels) will be provided with the clinical system. This will allow rapid viewing of the mammogram as a complete image at reduced resolution for adjustment of display parameters or with full 50 µm resolution in a region of interest which can be positioned with a trackball over any part of the image. The image output will also be provided as a laser-printed film image (4k x 5k pixels). The radiologist will be able to manipulate the display on the monitor to define the display characteristics of the image to be printed on the film.

The clinical version of the digital mammography system is still under construction. However, some preliminary measurements have been made on a prototype. The system can acquire a mammogram in 3-6 seconds with a dose to the breast of 0.85 mGy or less. Resolution has been measured at 9.5 line-pairs/mm. The dynamic range of the CCD is 5000:1 and with digitization of 12 bits a range of over 100 in x-ray exposure transmitted by the breast can be accommodated with a "worst case" display capability of over 40 shades of gray even in the densest part of the breast, depending on the level of quantum noise.

### 3. Computer Assisted Diagnosis (CAD)

CAD refers to a diagnosis made by a radiologist who takes into consideration the results of a computerized analysis of radiographic images and uses them as a "second opinion" [4,5]. The goal of CAD is to reduce the screening load and improve the diagnostic accuracy by reducing the number of false negative diagnoses. Preliminary results indicate that computers can aid in recognizing abnormalities and actually point out suspicious findings such as microcalcifications and masses. Several computer techniques have been applied to mammograms including: automatic (operator and image independent) enhancement methods for outlining

specific features such as normal parenchymal tissues, microcalcifications, and suspicious areas and pattern recognition and image segmentation methods for automatic localization, detection, and classification of suspicious breast lesions or normal parenchymal tissues.

In developing machine-assisted screening and diagnostic methods, one strives for automatic techniques which are both sensitive and specific, since the consequences of false negative interpretation (missed cancers) and false positive (FP) interpretation (traumatic, expensive investigation) are both serious. The goal of the computerized methods is to improve the performance of the radiologist by noise suppression, detail preservation, edge detection, and contrast enhancement and standardize the methods for image interpretation.

Automatic CAD schemes are the ultimate goal in mammography. Their development faces problems such as high-false positive detection rates, long CPU times, limited databases, low quality display devices. For such a workstation to be successful, one requires: high quality digital mammograms, high speed computers, large databases, efficient image processing techniques, characterization of image features of normal and abnormal patterns, understanding of image interpretation process by radiologists [4,5]. Intelligent radiologic workstations will not only retrieve, display, and process images but will also provide a wide range of tools to help us think more effectively about radiologic problems. This implies the inclusion of case-specific background information, reference images, consultations and new information from the literature.

It is anticipated that various methods with varying degrees of complexity will be required for optimum image enhancement, segmentation and pattern recognition of the mammographic features. The implementation of many of these algorithms will demand extensive computation times. Very large scale integrated (VLSI) circuits and image compression algorithms may provide a fast and cost effective technological solution to these computer vision and image processing areas. Some tasks accomplished until now are described in the following.

#### **A. Automated Detection of Clustered Microcalcifications**

At the University of Chicago (UC), a computer program is being developed to automatically locate clustered microcalcifications on mammograms [15,16,17,18]. With this method, a digital mammogram is processed by a linear filter to improve the signal-to-noise ratio of microcalcifications on the image. Gray-level thresholding techniques, which combine a global gray-level thresholding procedure and a locally adaptive gray-level thresholding procedure, are then employed to extract potential signal sites from the noise background. Subsequently, signal-extraction criteria are imposed on the potential signals to distinguish true signals from noise or artifacts. The computer then indicates locations that may contain clusters of microcalcifications on the image.

For 60 mammograms used in the study, the true-positive (TP) cluster detection accuracy of our automated detection program reached 85% at an FP detection rate of 2 clusters per image. An ROC study was performed to determine whether this performance level could result in an improvement in radiologists' performance when the CAD results were displayed on images. The results of the ROC study, as shown in Figure 2, indicated that CAD does significantly improve radiologists' accuracy in detecting clustered microcalcifications under conditions that simulate the rapid interpretation of screening mammograms. The results suggested also that a reduction in the computer's false-positive rate will further improve radiologist's diagnostic accuracy, although this improvement fell short of statistical significance.

At the Center for Engineering and Medical Image Analysis (CEMIA) at the University

of South Florida (USF), two-channel and three-channel quadrature mirror filters are developed for image decomposition and reconstruction [19,20] and dynamic neural networks are implemented for breast feature detection and extraction [21]. The sensitivity and specificity of detection is very high with these approaches. A preliminary study with 15 mammograms each containing at least one calcification cluster showed a TP rate of 100% with only 0.1 FP clusters per image; application of these methods to a larger data set is currently under way for a fuller evaluation.

## **B. Automated Detection of Mammographic Masses**

Similarly, a computerized scheme is under development in UC for the detection of masses in digital mammograms [22,23,24]. Based on the deviation from the normal architectural symmetry of the right and left breasts, a bilateral-subtraction technique is used to enhance the conspicuity of possible masses. The scheme employs two pairs of conventional screen-film mammograms (the right and left MLO views and CC views), which are digitized. After the right and left breast images in each pair are aligned, a nonlinear bilateral-subtraction technique is employed that involves linking multiple subtracted images to locate initial candidate masses. Various feature-extraction techniques are then used to reduce false-positive detections resulting from the bilateral subtraction. In an evaluation study using 154 pairs of clinical mammograms, the scheme yielded an 85% TP rate at an average of 3 false-positive detections per image.

Alternatively, tree-structured nonlinear filters, quasi-range edge detectors, and wavelets (two-channel quadrature mirror filters) are developed in CEMIA at USF and used for enhancement and edge detection of circumscribed, irregular and stellate masses. Preliminary results on a small number of mammograms show improved performances of these algorithms for noise suppression with simultaneous image detail preservation [25].

## **4. Telemammography**

Telemammography faces all the challenges associated with the acquisition, storage, transmission, processing and display of large amounts of data. The resolution and dynamic range currently required for digital representation of chest and musculoskeletal radiography (image size 2k x 2k to 4k x 4k pixels, pixel intensity encoded in 10 to 12 bits) stress both storage capacities and bandwidth capabilities of existing picture archiving and communications systems (PACS). Data compression studies and applications have been reported in the medical literature for over 10 years. Fidelity criteria are currently based on observer performance studies using selected case material and board certified radiologists as observers. Both lossless and lossy techniques have been offered in commercial systems [26]. Lossy techniques are generally used in cases not requiring primary interpretation from the lossy data set. Although standards are in the process of being developed (American College of Radiologists - National Electrical Manufacturers Association (ACR-NEMA) joint committee to develop a Standard for Digital Imaging and Communications in Medicine), the struggle between lossless and lossy compression techniques continues [27].

A screening mammography test consists of at least four images with each digital image ranging from 16 Mpixels to 64 Mpixels with dynamic ranges of 10 to 16 bits per pixel. Such large data sets and the fidelity requirements of mammography challenge the storage and bandwidth capabilities of existing communication systems. Cost effective storage of these images

and responsive image delivery via telecommunications channels can be facilitated using data compression technologies. Factors of 2 to 3 for lossless storage and transmission may be supported by existing encoder/decoder implementations and require no compromise in image fidelity. Significant cost savings in both research and clinical database storage is likely to result. Lossy compression approaches, offering higher gains, will need to be evaluated against observer performance for visual presentation and primary interpretation of mammographic data and evaluated for applications involving additional image processing and CAD.

Most of the results reported for lossless compression achieve a factor of 2 to 3 in compression ratio. A number of carefully constructed observer performance studies have reported successes with a "visually lossless" presentation of the data using lossy compression techniques. Block oriented discrete cosine transforms (16 x 16 to full frame) coupled with various adaptive encoding strategies produced compression results in the range of 20-to-30 :1 with no statistically significant differences in radiologist performance [28,29,30,31]. Compression ratios of 2-3 have been reported with lossless methods, e.g. tree-based codes.

Recent mammogram compression studies using wavelets also show promise with ratios up to 70:1 depending on the image [32]. Figures 3 and 4 present examples of mammogram compression with Haar wavelets at different rates. The original images (Figs. 3(a) and 4(a)) contain clustered microcalcifications and are compressed at 25:1 (Fig. 3(b)) and 50:1 rates (Fig. 4(b)). Although some image detail is lost with this type of wavelet compression, the appearance of the microcalcification clusters is not significantly affected. Furthermore, processing the original and the compressed reconstructed images with two channel wavelets results in similar segmentation of the calcification cluster from either image despite the losses during compression [25]. These results indicate that image processing of the compressed data could partially compensate for the information loss and encourage the acceptance of "visually lossless" compressed images.

## 5. Dual-Use Technologies

The results presented in the previous sections are representative of current preliminary work in digital mammography. It is anticipated that alternative approaches could be identified or developed which may be more successful and worthy of further study. In this context, a survey of technologies currently used in NASA Centers and Federal Research Laboratories was undertaken and has identified several projects that are promising and may have an impact in mammography. These projects span all the areas of interest and include: (a) scanning slot detectors using glass or plastic scintillating micro-fiber plates as the x-ray converting material and fiber-optic coupling to a CCD camera [33], silicon or amorphous-silicon arrays, and other advanced digital sensors for x-ray imaging; (b) software packages and algorithms such as neural networks, wavelets, and Bayesian classifiers used for target or object detection ; (c) lossless and lossy compression algorithms for handling large amounts of space image data, real-time software and systems for telemetry applications; (d) storage devices and local area networks to transmit real-time voice and video traffic with simultaneous transmission of computer data; (e) VLSI circuits suitable for implementing wavelets and neural networks for pattern recognition and compression problems in real time; and (f) telerobotic developments with potential applications to stereotactic mammography procedures. The idea of technology transfer is, therefore, realistic, and is expected to receive increasingly enthusiastic response.

## References

- [1] S. Shapiro, W. Venet, P. Strax, L. Venet: "Current results of the breast cancer screening randomized trial: the Health Insurance Plan (HIP) of Greater New York study" in *Screening for Breast Cancer*, N. E. Day and A. B. Miller eds., Toronto: Hogrefe International, 1988.
- [2] L. Tabar *et al*: "Reduction in mortality from breast cancer after mass screening with mammography," *Lancet*, vol. 12, pp. 829-832, 1985.
- [3] H. Seidman, *et al*: "Survival experience in the Breast Cancer Demonstration Project," *Cancer*, vol. 37, pp. 258-290, 1987.
- [4] F. Shtern: "Digital mammography and related technologies: a perspective from the National Cancer Institute," *Radiol.*, vol. 183, pp. 629-630, 1992.
- [5] K. Doi: "Computer-aided image interpretation," presented at *Breast Imaging: State-of-the-art and Technologies of the Future*, Bethesda, MD, September 4-6, 1991.
- [6] R. M. Nishikawa, M. J. Yaffe: "Signal-to-noise properties of mammography film-screen systems," *Med. Phys.*, vol. 12, pp. 32-39, 1985.
- [7] H. Kato: "Photostimulable phosphor radiography design considerations," in *Specification, acceptance testing and quality control of diagnostic x-ray imaging equipment*, J. A. Siebert, G. T. Barnes, and R. G. Gould, eds., Proc. 1991 Summer School, American Association of Physicists in Medicine (AAPM), in press.
- [8] J. Janesick, T. Elliot, A. Dingizian, *et al*: "New advancements in charge-coupled device technology - sub-electron noise and 4096x4096 pixel CCDs," Proc. SPIE Symp. on Electronic Imaging, vol. 1242, 1990.
- [9] R. A. Street, *Hydrogenated amorphous silicon*, Cambridge: Cambridge University Press (ISBN 0 521 37156 2), pp. 391, 1991.
- [10] R. A. Street, S. Nelson, L. Antonuk, V. Perez-Mendez: "Amorphous silicon sensor arrays for radiation imaging," Proc. MRS Symp., 1990.
- [11] J. A. Rowlands, G. DeCrescenzo, N. Araj: "X-ray imaging using amorphous selenium: Determination of x-ray sensitivity by pulse height spectroscopy," *Med. Phys.*, vol. 19(4), pp. 1065-1069, 1992.
- [12] R. S. Nelson, Z. Barbaric, L. W. Bassett, R. Zach: "Digital slot scan mammography using CCDs," Proc. SPIE, vol. 767, pp. 102-108, 1987.
- [13] M.J. Yaffe: "Digital Mammography," in *Syllabus of Categorical Course on Technical Aspects of Mammography*, A. Haus and M. J. Yaffe, eds., Oak Brook, IL: Radiological Society of North America, 1992.
- [14] A. D. A. Maidment, D.B. Plewes, B.G. Starkoski, G.E. Mawdsley, I.C. Soutar and M.J. Yaffe: "A Time Delay Integration Imaging System for Digital Mammography," *Med. Phys.*, 1992 (submitted for publication).
- [15] H. P. Chan, K. Doi, S. Galhotra, C. J. Vyborny, H. MacMahon, P. M. Jokich: "Image feature analysis and computer-aided diagnosis in digital radiography: 1. Automated detection of microcalcifications in mammography," *Med. Phys.*, vol. 14, pp. 538-548, 1987.
- [16] H. P. Chan, K. Doi, C. J. Vyborny, K. L. Lam, R. A. Schmidt: "Computer-aided detection of microcalcifications in mammograms: methodology and preliminary clinical study," *Invest. Radiol.*, vol. 23, pp. 664-671, 1988.
- [17] H. P. Chan, K. Doi, C. J. Vyborny, R. A. Schmidt, C. E. Metz, K. L. Lam, T. Ogura, Y.

- Wu, H. MacMahon: "Improvement in radiologists' detection of clustered microcalcifications on mammograms: the potential of computer-aided diagnosis," *Invest. Radiol.*, vol. 25, pp. 1102-1110, 1990.
- [18] R. M. Nishikawa, Y. Jaing, M. L. Giger, K. Doi, C. J. Vyborny, R. A. Schmidt: "Computer-aided detection of clustered microcalcifications," *Proc. IEEE ICSMC-92*, pp. 1375-1378, 1992.
- [19] W. Qian, L. P. Clarke, M. Kallergi, H-D. Li, R. P. Velthuizen, and R. A. Clark: "Tree-structured nonlinear filter and wavelet transform for microcalcification segmentation in mammography," *Proc. IS&T/SPIE Annual Symposium on Electronic Imaging, Science & Technology*, San Jose, California; January 31 - February 5, 1993.
- [20] W. Qian, L. P. Clarke, H-D. Li, M. Kallergi, R. P. Velthuizen, R. A. Clark, and M. L. Silbiger: "Digital mammography: tree-structured wavelet decomposition and reconstruction for feature extraction," *Int. J. Pat. Rec. Art. Intel.*, 1993 (submitted for publication).
- [21] K. S. Woods, J. L. Solka, C. E. Priebe, C. C. Doss, K. W. Bowyer, and L. P. Clarke: "Comparative evaluation of pattern recognition techniques for detection of microcalcifications," *Proc. IS&T/SPIE Annual Symposium on Electronic Imaging, Science & Technology*, San Jose, California; January 31 - February 5, 1993.
- [22] M. L. Giger, F-F Yin, K. Doi, C. E. Metz, R. A. Schmidt, C. J. Vyborny: "Investigation of methods for the computerized detection and analysis of mammographic masses," *Proc. SPIE*, vol. 1233, pp. 183-184, 1990.
- [23] M. L. Giger, R. M. Nishikawa, K. Doi, F-F Yin, C. J. Vyborny, R. A. Schmidt, C. E. Metz, Y. Wu, H. MacMahon, H. Yoshimura: "Development of a "smart" workstation for use in mammography," *Proc. SPIE*, vol. 1445, pp. 101-103, 1991.
- [24] F-F Yin, M. L. Giger, K. Doi, C. E. Metz, C. J. Vyborny, R. A. Schmidt: "Computerized detection of masses in digital mammograms: analysis of bilateral subtraction images," *Med. Phys.*, vol. 18, pp. 955-963, 1991.
- [25] L. P. Clarke: "Digital mammography: advances in image restoration, enhancement, and feature extraction using NN's, wavelets, and fuzzy logic approaches," *President's Symposium - New developments in electronic imaging: Mammography*, *Proc. AAPM 35th Annual Meeting*, August 8-12, Washington, DC, 1993.
- [26] D. L. Wilson: "Compression for radiological images," *Proc. SPIE Conf. on Medical Imaging VI: PACS Design and Evaluation*, vol. 1654, pp. 130-139, 1992.
- [27] "ACR-NEMA Digital Imaging and Communication Standard Committee, Working Group #4, MEdPACS section," *Data Compression Standard #PS2*, Hartwig Blume, Chairperson, 1989.
- [28] H. MacMahon, K. Doi, S. Sanada, S. M. Montner, M. L. Giger, C. E. Metz, N. Nakamori, F-F, Yin, X. W. Xu, H. Yonekawa, H. Takeuchi: "Data compression: effect of diagnostic accuracy in digital chest radiography," *Radiol.*, vol. 178, pp. 175-179, 1991.
- [29] T. Ishigaki, S. Sakuma, M. Ikeda, Y. Itoh, M. Suzuki, S. Iwai: "Clinical evaluation of irreversible image compression: analysis of chest imaging with computed radiography," *Radiol.*, vol. 175, pp. 739-743, 1990.
- [30] J. Sayre, D. R. Aberle, M. I. Boechat, T. R. Hall, H. K. Huang, B. K. Ho, P. Kashfian, G. Rahbar: "Effect of data compression of diagnostic accuracy in digital hand and chest radiography," *Proc. SPIE Conf. on Image Capture, Formatting and Display*, vol. 1653, pp. 232-240, 1992.
- [31] J. Chen, M. J. Flynn: "The effect of block size on image quality for compressed chest

- radiographs," Proc. SPIE Conf. on Image Capture, Formatting and Display, vol. 1653, pp. 252-260, 1992.
- [32] R. DeVore, B. Jawerth and B. Lucier: "Image compression through wavelet transform coding," IEEE Trans. Inf. Theory, vol. 38, pp. 719-746, 1992.
- [33] J. K. Walker: "Scintillation fiber technology for high resolution digital diagnostic x-ray applications," President's Symposium - New developments in electronic imaging: Mammography, Proc. AAPM 35th Annual Meeting, August 8-12, Washington, DC, 1993.

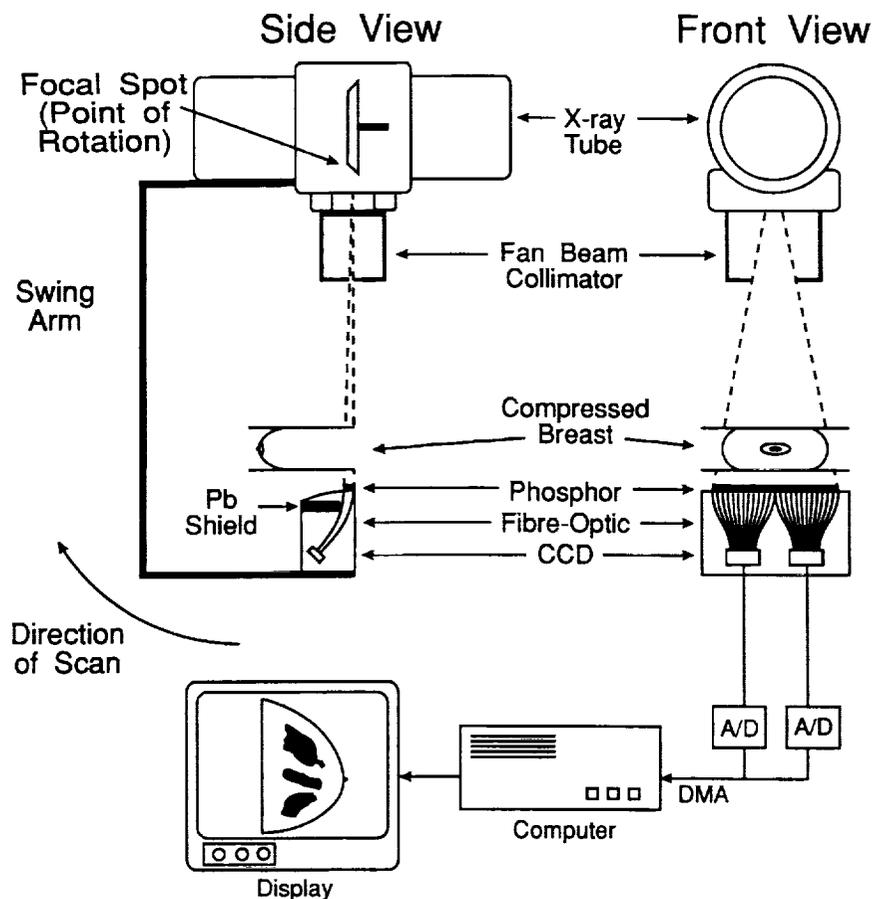


Figure 1. Schematic diagram of scanned-slot digital mammography system.

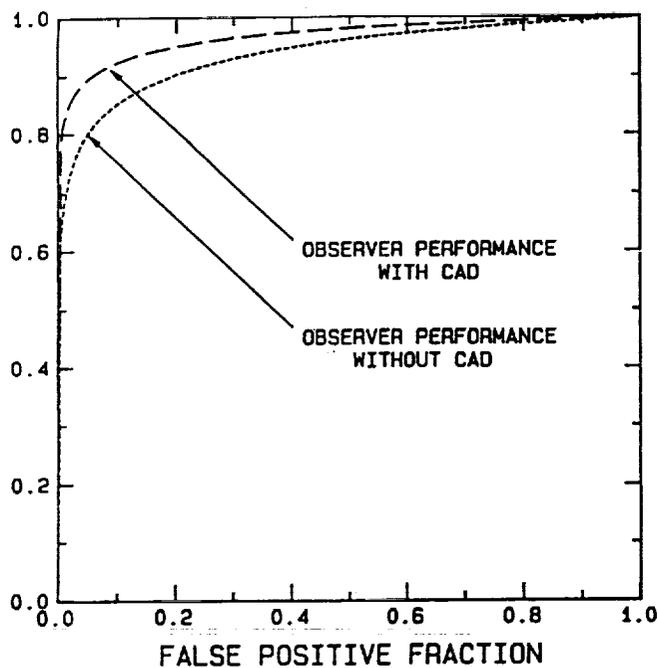
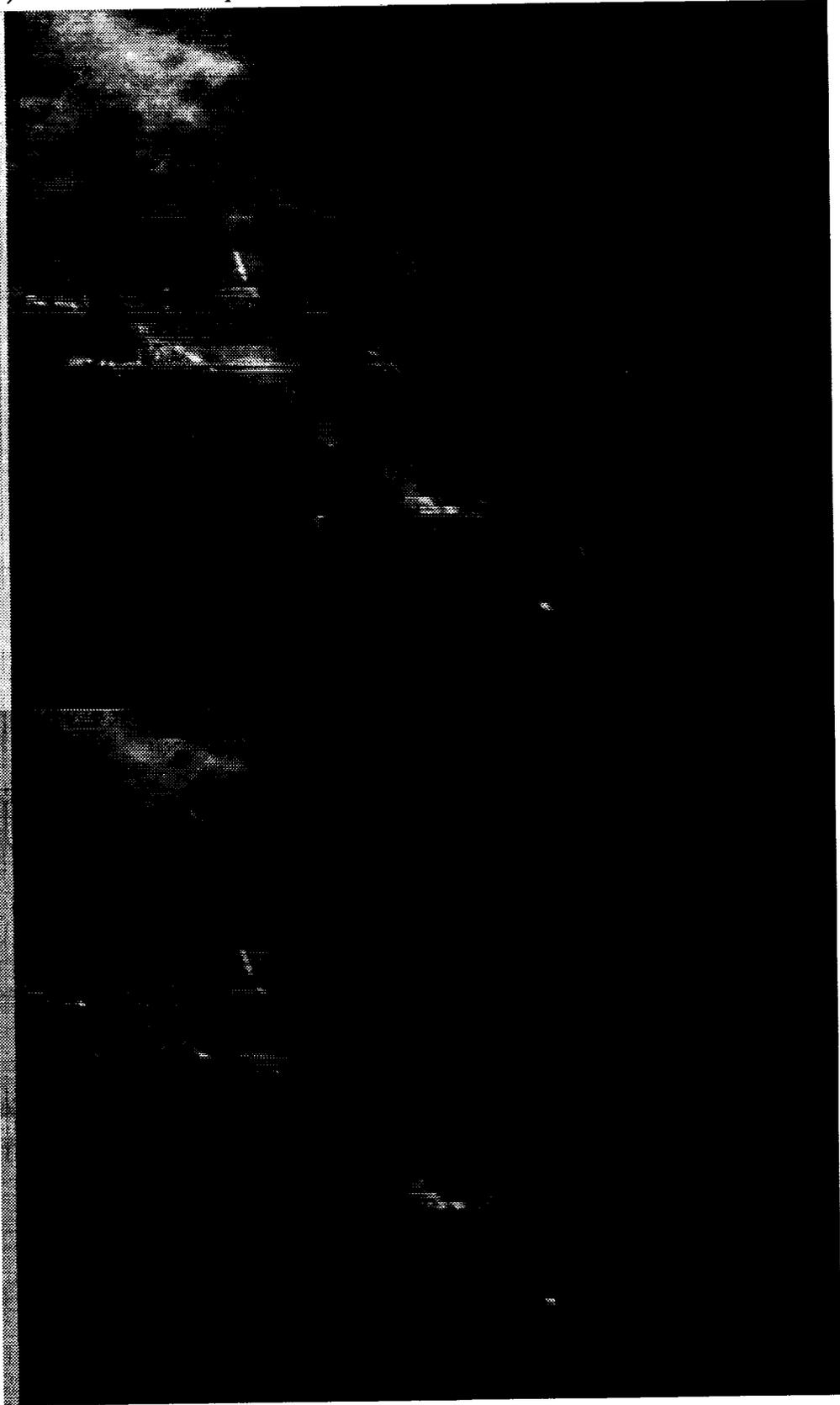


Figure 2. Comparison of ROC curves for two reading conditions (with and without computer output) for observer performance study in detecting microcalcification clusters in mammograms.

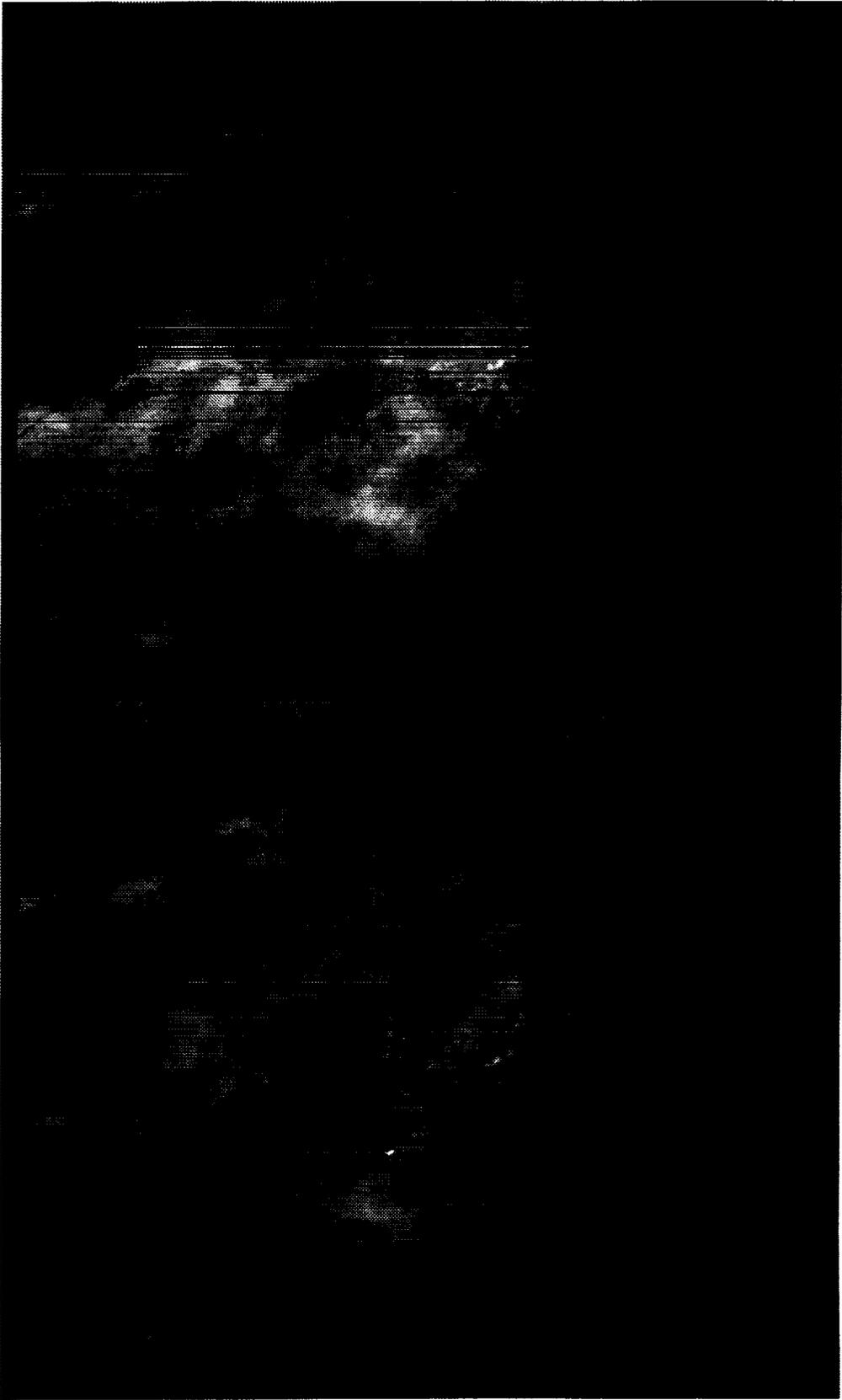
**Figure 3.** (a) Section of original digitized mammogram with clustered microcalcifications and (b) Haar wavelet compressed and reconstructed image at a rate of 25:1



(a)

(b)

**Figure 4.** (a) Section of original digitized mammogram with clustered microcalcifications and (b) Haar wavelet compressed and reconstructed image at a rate of 50:1



(a)

(b)

## A STUDY OF VIDEO FRAME RATE ON THE PERCEPTION OF MOVING IMAGERY DETAIL

Richard F. Haines  
RECOM Technologies  
Ames Research Center - NASA  
Moffett Field, California 94035

Sherry L. Chuang  
Spacecraft Data Systems Research Branch  
Ames Research Center - NASA  
Moffett Field, California 94035

### Abstract

The rate at which each frame of color moving video imagery is displayed was varied in small steps to determine what is the minimal acceptable frame rate for life scientists viewing white rats within a small enclosure. Two, twenty five second-long scenes (slow and fast animal motions) were evaluated by nine NASA principal investigators and animal care technicians. The mean minimum acceptable frame rate across these subjects was 3.9 fps both for the slow and fast moving animal scenes. The highest single trial frame rate averaged across all subjects for the slow and the fast scene was 6.2 and 4.8, respectively. Further research is called for in which frame rate, image size, and color/gray scale depth are covaried during the same observation period.

### Introduction

The perception of moving detail(s) on a computer monitor or TV screen is a complex function of many optical, visual, and cognitive variables; disagreement remains concerning the impact of specific variables. For example Farrell and Booth (1984) reported that decreasing video bandwidth produces relatively little reduction in subjectively determined image acuity for moving objects while Connor and Berrang's (1974) data suggest a linear relationship between increased bandwidth and increased judged image quality. Some investigators feel that this linear relationship results from an improvement in perceptibility due to increasing speed of image motion across the screen. However, given the same amount of bandwidth reduction and speed of image motion, the impairment of image quality is greater for images having many vertically oriented edges of high contrast than for images with only a few such edges. So both the contrast and orientation of the objects are important.

Initially we assumed that those who work with small animals prefer to see smoothly moving images rather than disjointed, choppy motion since smooth motion supports improved image recognition and more correct interpretation of behavioral functions and interactions.

A number of other earlier studies have been performed on the effect of varying frame rate on image usefulness. Ranadive (1979) reported that video bandwidth was directly proportional to the product of *resolution* (height x width; pixels per frame), *frame rate* (fps), and *gray scale* (bits/pixel). When the viewer varied one of these three parameters at a time (while watching his own motions controlling a robot in order to perform a simple task), it was found that he could carry out the assigned task relatively well even though these image parameters were degraded significantly. Performance was defined as the quotient  $T_t/T_d$  where  $T_t$  is the time to accomplish the task using full video (i.e., no degradation) and  $T_d$  is the time required to accomplish the task using degraded video. He found that when only one of the three parameters was systematically reduced performance

remained at acceptable levels until a point was reached where the task could no longer be accomplished at all. He also found that frame rate and gray scale could be degraded by larger amounts than resolution before the critical performance limit was reached. Since the total bits associated with the frame rate parameters in Ranadive's study was only 42 percent of the total bits associated with the other two parameters this suggests that frame rate is a very attractive candidate for reducing video bandwidth under these viewing conditions.

Deghuee (1980) had an operator adjust resolution, frame rate, and gray scale during manual robotic control operations under total bit rate constraints. Dynamically changing these three parameters in real time influenced performance *although lower bit rates did not result in reduced performance*. Since only two bit rates were studied (10 kbps and 20 kbps) it is possible that these total bit rate conditions were not sufficiently small enough and/or sufficiently different from one another to produce significant decrements in performance. Deghuee also reported that the operators did not adjust the three parameters to achieve an image with some "optimal" quality but, rather, set each parameter to achieve some predetermined combination of settings of the three available parameters. Because his operators were sufficiently familiar with the appearance of changes in each of the three parameters separately they were (probably) able to adequately anticipate the appearance of a predetermined combination of them. Deghuee also found that the type of manipulation task undertaken yielded the most significant differences in performance which is what we found when comparing different levels of video compression (Haines and Chuang, 1992).

None of the studies cited above varied frame rate systematically while viewers evaluated the health and behavior of small animals as will be done in future Space Station Freedom experiments. This paper describes a study of the relationship between video frame rate and perceived quality and acceptability to life scientists of moving imagery of white rats. It is another in a continuing series of studies related to remote monitoring between earth orbit and the ground where transmission bandwidth is limited and must be used optimally.

As Haskell and Steele (1981) state, "Only when perception is properly understood will we have accurate objective measures. However, the day when we can, with confidence, objectively evaluate a new impairment without recourse to subjective testing seems very remote." The interested reader should consult (Gonzalez and Wintz, 1987; Watson, 1987; Watson et al., 1983; Wood et al., 1971) for further information on this issue.

## Method

*Experimental Design and Variables.* The experimental design used may be characterized as a 2 x 3 x 2 x 9 parametric design having the following factors:

- 2 levels of direction of change of frame rate (increasing; decreasing)
- 3 levels of frame rate change resolution (5, 2, 1 fps)
- 2 scenes (slow animal motion; fast animal motion)
- 9 subjects (Ss)

Each subject (S) was presented all twelve cell conditions. Five subjects received scene 1 first while the other four received scene 2 first. Likewise, four subjects received increasing frame rate trials first per pair while the other five received decreasing frame rate trials first. Frame rates from 1.5 to 30 fps were explored.

The method of limits (Woodworth and Schlosberg, 1965) was used to quantify the effect of video frame rate on perceived image quality. This method employs alternating series of

decreasing and increasing frame rates where S indicated the frame rate at which he or she could no longer accept the quality of the moving imagery and then gave a numeric rating of image quality at each frame rate presented. Each series of trials was conducted at progressively smaller frame rate steps: Initial trials varied in five fps steps in order to quickly identify the approximate frame rate separating an acceptable from an unacceptable image. Subsequent trials varied in 2 fps and 1 fps steps. Thus, S was progressively exposed to finer and finer frame rate steps. Means of the 2 fps and 1 fps trials were combined to determine the final threshold frame rate for each subject. Two separate judgments were made immediately following each 25 second-long scene finished:

- (1) Was the scene of acceptable quality to make useful scientific judgments in their own scientific discipline (yes, no)?
- (2) What was the image quality? A five point scale of whole numbers was used: (1) = image clarity completely unacceptable relative to 30 fps, (3) = image clarity is of average acceptability relative to 30 fps, and (5) image clarity is completely acceptable relative to 30 fps.

*Video Tape Scene Description.* The so-called "slow scene" showed two white rats within a small enclosure. Almost all of the scene showed the animals performing typical grooming activity (e.g., licking their fur, scratching with a hind leg at about 6 - 10 Hz, playfully biting each other). Neither animal walked around very much during the scene but exhibited typical slow limb and body movement, exploratory behavior such as sniffing, etc. The so-called "fast scene" showed the same white rats inside the same enclosure but they were engaged in playful behavior such as tumbling, chasing and rolling over each other, and mock fighting during most of the scene. The angular rates of some of their movements were so great that they appeared to be almost at the edge of blurring, viewed at 30 fps.

*Procedure.* A training and familiarization period was provided where the scene to be evaluated was presented many times (typically five to seven) on an 18" color standard television monitor at 30 fps so that the subject could become very familiar with it. An experimenter discussed the objective of the study and answered questions during this time. The subject was also asked to write down what scene details were of importance and which would be used to evaluate the scene. The objective was to try to ensure that the same scene-judgement criteria would be used throughout the study. This objective was also emphasized verbally prior to data collection.

A decreasing frame rate test run began with a twenty five second-long scene at 30 fps followed by another identical twenty five second-long scene at 25 fps, etc. Judgements were made immediately following each scene presentation. This procedure continued until the subject indicated that the scene details were no longer acceptable to them to make useful scientific judgments in their scientific discipline. This was followed immediately by an ascending series of trials beginning with the smallest frame rate. A ten second-long period of gray screen occurred between each scene presentation during which S looked away from the screen and verbalized his or her ratings and the experimenter changed the conditions for the next trial and recorded S's ratings. Another increasing and decreasing series of trials followed immediately in which frame rate was varied in 2 fps steps. A final series of increasing and decreasing trials then followed in 1 fps steps. The starting fps for the 2 and 1 fps step trials were estimated on the basis of each S's judgments made during the earlier trials.

*Subjects.* Nine volunteers took place, 5 male (minimum = 38 yrs; maximum = 56 yrs;

mean age = 50) and 4 female (minimum = 28 yrs; maximum = 42 yrs; mean age = 33.5). All possessed 20:20 corrected or uncorrected distance acuity and normal color perception. Two had taken part in previous video compression studies conducted by the authors.

*Apparatus.* All imagery was presented on a 16" (diagonal) VGA screen of the IBM computer. This PS/2 Model 80-321 computer has 10 megabytes (MB) of RAM and a 320 MB hard disk. The video imaging hardware installed in it consisted of Intel's "ActionMedia II" board set; an Action Media II Capture module attaches to the ActionMedia II Delivery Board as a daughter board. (FN-1) The prerecorded analog video segments (scenes) described above were played on a four-head, Heliquad II Model JR4500 VHS video cassette recorder whose video output was connected to the composite RS170 input connector of the ActionMedia II boards. They were displayed in a small inset video window measured 5.25" (h) by 3.75" (w) on the larger computer monitor and subtended 12.5 degrees horizontally and 9 degrees vertically (of the observer's visual field).

A software application by IBM known as "Person-to-Person" was used in conjunction with the digital imaging hardware. This application runs with OS/2's Presentation Manager and permits live video to be displayed within an on-screen video window in the video conferencing mode. The following video settings were used: Tint = 50%, Saturation = 76%, Brightness = 66%, and Contrast = 50%, View = single, Effects = local, Large View. An on-screen frame rate control was used which allowed a frame rate to be selected between 30 frames per second and 1.5 frames per second.

All video imagery was compressed using a nine bit hardware-based compression technology developed jointly by IBM and Intel Corporation known as Digital Video Interactive (DVI). This compression approach divides each video frame into four by four pixel blocks and allocated one pixel representation. The pixel representation consists of eight bits for luminance and one bit for hue (color) and saturation. This algorithm is used within each frame i.e., no interframe encoding. Because the scenes presented here were repeated, identical twenty five second-long segments, the only perceptually relevant parameter that changed from trial to trial was frame rate.

## Results

The results are presented in three sections: I. Mean image acceptance results, II. Highest Frame rate at which image quality was totally unacceptable, and III. Image evaluation criteria used.

*I. Mean Image Acceptance Results.* Table 1 presents the minimum acceptable frame rate (averaged across all trials per S) for each type of scene. Experience category, age and sex are also given for each S. The raw data are given in Appendix A and B. It can be seen that: (1) these Ss accepted image quality at frame rates between 1.5 to 8.5 fps. Indeed, the three most highly experienced Ss felt that they could obtain all needed information at rates *below* 1.5 fps which was the slowest rate possible from our hardware. (2) the slow versus fast animal scene did not yield a statistically significant difference in acceptable mean minimal frame rate across all Ss. However, four of the Ss did require a faster frame rate for the fast scene of about one fps, (3) when these data were grouped by general level of familiarity and experience with white rats, mean acceptance frame rate was not clearly different either for the slow or the fast scene across these experience levels, and (4) there was no significant difference between the male and female S's mean data.

Table 1

Mean Minimum Image Acceptance Results (fps)  
for Each Subject Averaged Across 2 fps and 1 fps Trials

Experience Category (note 1)	Subj. No.	Age	Sex	Slow	Fast	
A	7	45	M	<1.5 (2)	<1.5	
A	1	55	M	<1.5	<1.5	
A	8	56	M	<1.5	<1.5	
B	4	28	F	4.9	6.0	
B	2	34	F	3.0	4.3	
B	5	38	M	3.9	4.9	
B	9	42	F	8.5	6.4	
B	3	56	M	5.6	3.7	
C	6	30	F	5.0	5.1	
				Mean =	3.9 (3)	3.9 (3)
				SD =	2.4	2.0

Footnotes:

1. A = 15 or more years of experience; B = 5 - 15 years; C = 0 - 5 years.
2. All values labelled < 1.5 were scored as 1.5.
3. Not statistically significantly different (t test).

*II. Highest Frame rate at which image quality was totally unacceptable.* This numeric rating provided a second response measure of the subjective usefulness or non-usefulness of low video frame rates. We are mainly concerned with the single highest frame rate that was judged to be of completely unacceptable image clarity. Table 2 and 3 provides these data.

Table 2

Highest Frame Rate Single Trial Judged to Provide a Totally  
Unacceptable Image Quality for the *Slow* Scene  
(Relative to 30 fps)

Subj. No.	Ascending Trials	Descending Trials
1	5.1	5.5
2	6.4	7.5
3	4.2	3.1
4	*	*
5	5.1	10.2
6	*	*
7	*	*
8	13.5	3.6
9	3.6	5.2
Mean =		6.3
		6.0
Grand Mean = 6.2		

\* Indicates that subject's fastest unacceptable frame rate was <1.5.

In addition to the above results it was found that: (a) there were characteristic individual differences in these numeric ratings. Each S gave consistent numeric ratings throughout their viewing period and did not appear to change their judgment criteria. This was shown by the fact that the same numeric score tended to be assigned to the same frame rate over time even though they had viewed different frame rates in the meantime, and (b) the Ss appeared to have understood and followed these rating instructions.

The grand mean data of Table 2 and 3 reinforce the previous Table 1 data with regard to the frame rate - scene motion relationship, viz., the slower scenes required a higher frame rate in order to be judged as acceptable by these Ss.

Table 3  
Highest Frame Rate Single Trial Judged to Provide a Totally  
Unacceptable Image Quality for the *Fast* Scene  
(Relative to 30 fps)

Subj. No.	Ascending Trials	Descending Trials
1	3.6	6.4
2	4.2	5.5
3	3.1	3.1
4	*	*
5	*	8.5
6	*	*
7	*	*
8	*	*
9	*	*
Mean = 3.6		5.9
Grand Mean = 4.8		

Footnotes:

\* Indicates that subject's fastest unacceptable frame rate was <1.5.

III. *Image Evaluation Criteria Used* (Professional Discipline, Experience Level, and Minimal Frame Rate). It was expected that each subject might use a somewhat different set of criteria for evaluating the moving imagery of each scene. Such differences might reflect differences in one's disciplinary training and professional experience. This was found to be the case. In fact, large individual differences were found in the minimum acceptable frame rate people selected during their scene evaluations. Having a lot of prior experience seemed to play an important role in making these judgements, perhaps by improving one's capability to extract subtle image cues or ignoring distracting cues that are present. For example, the three Ss who possessed the most research experience also had prior experience in viewing one (1) fps images of rats in micro-gravity. They judged all scenes at 1.5 fps and higher as being entirely adequate for making their judgments of grooming

behavior, general weight and health of the animals, evidence of edema and porphorin (exudate) build-up around the nose, ears and eyes, reaction to allergies, fecal matter build-up around the tail, and leg extension movements. Apparently, their prior experience permitted them to notice these details regardless of how quickly and discontinuously the image shifted across the screen. However, it must be noted that this particular list of image characteristics is made up mostly of static cues. Less experienced subjects generally required higher frame rates to make their judgments. This finding argues in favor of allowing each user to set his or her own frame rate, if possible, to support their own scientific requirements.

## Discussion

A minimum frame rate was identified in the present study where experienced subjects judged the quality of image motion (and other details) as being acceptable to them to adequately judge the overall status and behavior of rats. The minimal frame rate averaged across all subjects was approximately four fps for both the slow and fast scene. Minimal acceptable frame rates varied from 1.5 fps to 5.1 fps for both the slow and the fast scene. It is clear from this study that what is an acceptable minimal frame rate is directly related to at least three complex factors: (1) the type of visual discriminations that must be made from the frames, (2) the nature of the moving images to be examined, and (3) the level of experience one has in making the judgments. These visual-cognitive discriminations range from being very general (e.g., is the animal alive?) to highly specific (e.g., is the animal displaying specific signs of allergic reactions or vestibular dysfunction?).

More than one third of all of the judging criteria cited by these Ss were static in nature (e.g., nasal discharge, hair texture, signs of blood, posture). It is possible that the presentation of multiple frames per second actually impeded visual judgments of these specific kinds of image features. Thus, there is probably a class of static image details of importance to the S, a class of dynamic image details of importance to the S, and a third class in which both are relevant in varying degrees. This possibility suggests the need for further experimentation in which various mixes of cues ranging from static only to dynamic only be presented at different frame rates to see if it is possible to identify minimal frame rates within each class of image details.

*Visual Integration of Object Motion.* The perception of a moving image on a TV screen is actually the result of visually smoothing a series of time sampled (strobed) still image frames into an apparently continuous movement. As individual picture elements (pixels) making up the full frame each change in intensity and color the eye attempts to integrate them and to identify the meaning of this constantly changing array of luminous dots.

Image details may or may not appear to move across the screen depending upon many variables. For instance, the combination of visual angle and duration over which adjacently illuminated pixels appear to change determines whether the image is seen as a strobed (jumping) or continuously (smooth) moving image. Watson et al. (1983) has found that image sampling frequency (Hz) increases almost linearly with an increase in the angular velocity of an image seen on a screen in order to produce smooth motion rather than strobe motion. Images translating at about one degree arc per second must be sampled at about 30 Hz in order to appear to be moving smoothly.

It is interesting to speculate whether minimal acceptable frame rate may be related somehow to the time required for the visual system to extract information from a scene during a single glance. For instance, Senders et al. (1964) reported that the mean visual dwell time (FN-5)

for visual informational displays having information bandwidths from 0.05 to 0.48 Hz was 0.4 sec. Interestingly, several other studies of eye fixation dwell time on displays also have shown a mean duration of about 0.4 second across a wide range of display bandwidths (Harris and Christhif, 1980; Carbonell et al., 1968).

*Subject Variables.* There is little doubt that the human visual system is remarkably adept at extracting useful information from relatively degraded video imagery. If resolution is degraded, for example, perception probably shifts to lower spatial frequencies which incorporate slightly higher visual contrasts in order to perceive image translation across the scene.

*Application of Data to Space Station Freedom Operations.* The planned video downlink rate capacity for Space Station Freedom will be variable in the following five steps (Corder, 1992):

60 (full frame) fields per second	41.1 MB/s
30 (1/2 frame) fields per second	20.8 MB/s
15 (1/2 frame) fields per second	10.5 MB/s
7.5 (1/2 frame) fields per second	5.3 MB/s
1.875 (1/2 frame) fields per second	1.5 MB/s

Assuming a full frame video image format of 500 x 400 x 8 bits and 30 fps the required data rate would be 6 MB/s. Even without digital image compression, use of 7.5 fps (which is a higher frame rate than almost all of the present Ss accepted) would reduce the downlink data rate by a factor of 4 relative to the 30 fields per second data rate given here. If the Tracking/Data Relay Satellite's (TDRSS) Ku band maximum downlink rate is 43 Mb/s (5.37 MB/s) then without video compression it would support only one (NTSC) video channel. Clearly, the downlink bandwidth of all channels must be reduced significantly in order to be able to support all of the required control and monitoring functions planned. Reducing frame rate appears to be an acceptable means of accomplishing this objective in some research situations.

### Conclusions

We conclude from these findings that video bandwidth may be reduced from SSF to the ground by a factor of more than 4 times the normal 30 fields per second (approx. 4 fps) and still provide an acceptable image to the majority of scientists and animal care personnel. Observer prior experience plays a central role in determining minimal acceptable frame rate. It is not yet clear whether these data can be extrapolated to other life science animal specimens.

### References

- Carbonell, J.R., J.L. Ward, and J.W. Senders, 1968. A queuing model of visual sampling: Experimental validation. *IEEE Transactions on Man-Machine Systems*, MMS-9, Pp. 82-87.
- Connor, D.J., and J.E. Berrang, 1974. Resolution loss in video images. *NTC 74 Record*, (IEEE Publ. 74, CHO 902-7 CSCB, Pp. 54-60), San Diego, Calif.: Institute of Electrical and Electronics Engineers.

- Corder, E.L., 1992. Internal video subsystem overview. Presentation given at Payload Data Services Workshop, Huntsville, Alabama, August 5-6.
- Deghuee, B.J., 1980. Operator-adjustable frame rate, resolution, and gray scale trade-off in fixed-bandwidth remote manipulator control. Mass. Inst. of Technol., *M.S. Thesis*, (Department of Aeronautics and Astronautics), Boston, Massachusetts.
- Farrell, R.J., and J.M. Booth, 1984. *Design handbook for imagery interpretation equipment*. Boeing Aerospace Co, Document D180-19063-1, Seattle, Washington.
- Gonzalez, R.C., and P. Wintz, 1987. *Digital Image Processing*. 2nd Ed., Addison - Wesley Publ. Co., Menlo Park, Calif.
- Haines, R.F., and S.L. Chuang, 1992. The effects of video compression on acceptability of images for monitoring life sciences' experiments. *NASA Technical Paper 3239*.
- Harris, R.L., and D.M. Christhilf, 1980. What do pilots see in displays? *Proceedings of the Human Factors Society, 24th Annual Meeting*, Los Angeles, CA, Human Factors Society, Pp. 22-26.
- Haskell, B.G., and R. Steele, 1981. Audio and video bit-rate reduction. *Proceedings of the IEEE*, Vol. 69, No. 2, Pp. 252-26.
- Ranadive, V., 1979. Video resolution, frame rate and gray scale tradeoffs under limited bandwidth for undersea teleoperation, Mass. Inst. of Technol., *M.S. Thesis* (Department of Aeronautics and Astronautics), Boston, Massachusetts.
- Senders, J.W., J.E. Elkind, M.C. Grignetti, and R.P. Smallwood, 1964. An investigation of the visual sampling behavior of human observers. *NASA CR-434*, Bolt, Beranek & Newman, Cambridge, Mass.
- Watson, A.B., 1987. Efficiency of a model human image code. *J. Optical Society of America*, Series A, Vol. 4, No. 12, Pp. 2401-2417.
- Watson, A.B., A. Ahumada, Jr., and J.E. Farrell, 1983. The window of visibility: A psychophysical theory of fidelity in time-sampled visual motion displays. *NASA Technical Paper 2211*.
- Wood, C.B.B., J.R. Sanders, and D.T. Wright, 1971. Image unsteadiness in 16mm film for television. *Journal of the Society of Motion Picture and Television Engineers*, Vol. 80, Pp. 812-818.
- Woodworth, R.S., and H. Schlosberg, 1965. *Experimental Psychology*. Holt, Reinhart and Winston, Inc., New York.

#### Footnotes

1. One of the present test subjects served as an investigator on the SL-3 project and had a great deal of experience viewing 1 fps scenes.
2. ActionMedia II boards digitize and compress a video signal for display on a monitor

and/or storage on a hard disk. The boards used here employed a dual-chip, B-series i750 Video Display Processor.

3. Integration here refers to performing content associations and storing this information in visual memory.
4. The Nyquist theorem states that it is necessary and sufficient to visually sample signal at two times its bandwidth.
5. Visual dwell time refers to the duration over which no eye movement occurs.

**Systems Aspects of COBE Science Data Compression**

I. Freedman, E. Boggess, E. Seiler (Hughes STX).  
Cosmology Data Analysis Center  
Greenbelt, MD 20771

**Abstract.** A general approach to compression of diverse data from large scientific projects has been developed and this paper addresses the appropriate system and scientific constraints together with the algorithm development and test strategy. This framework has been implemented for the COsmic Background Explorer spacecraft (COBE) by retrofitting the existing VAX-based data management system with high-performance compression software permitting random access to the data.

Algorithms which incorporate scientific knowledge and consume relatively few system resources are preferred over *ad hoc* methods. COBE exceeded its planned storage by a large and growing factor and the retrieval of data significantly affects the processing, delaying the availability of data for scientific usage and software test. Embedded compression software is planned to make the project tractable by reducing the data storage volume to an acceptable level during normal processing.

**1. Introduction**

Large scientific projects generate diverse scientific, engineering and instrument housekeeping data at rates that frequently exceed the capacity of storage and retrieval devices. Although many techniques have been proposed in the data compression literature [1], almost all are based on data models that make predictions based on a few successive pixels or a few hundred images in a training set. These data models do not incorporate the *a-priori* scientific knowledge of approximate relations between data set elements (physical laws) or the known accuracy requirements for specific elements of record structures. Such knowledge reduces the specific entropy of the data, enabling an effective trade-off in wall-clock processing time between additional cycles for on-the fly compression and decompression and a reduced input-output load.

If the system response is sensitive to the network load (when the network is saturated) reduction in storage complexity may be as critical as reduction of the overall load. Furthermore, fixed mechanical disks are an expensive resource and the risk of catastrophic data loss increases dramatically with the number of disks on the system. Local SCSI disks are sometimes suggested to represent inexpensive storage media but the access time is relatively long. Mass storage devices such as magnetic tape juke boxes can be less than ideal as the tape quickly stretches with use and becomes unreadable after a short time (1 year) compared to the typical project lifetime (20 years) necessitating frequent and expensive data migration.

## 2. COBE Science Goals and Achievements.

The COsmic Background Explorer (COBE), NASA's first satellite devoted to the study of cosmology was launched on 18 November 1989. The cryogenic period of the mission covered the time from 21 November 1989 to 21 September 1990. COBE carries three instruments: the infrared experiment DIRBE, the anisotropy experiment DMR and the spectrum experiment FIRAS, of which DIRBE and DMR are still operating [2].

All three instruments have achieved their preliminary goals. FIRAS has shown the far infrared background to be isotropic to 0.03% and consistent with a black-body radiating at 2.726 K. [3] DMR has revealed further evidence of the Big Bang theory of cosmology in the form of a spectrum of ripples at the level of 1 part per million after known astrophysical foreground sources have been subtracted from the integrated signal. DIRBE has placed upper limits on the spectrum of the diffuse celestial background which are more stringent than previously available [4]. DIRBE and FIRAS have contributed to Galactic astronomy by mapping the stars in the direction of the Galactic core [5], modelling the physical conditions in the interstellar medium [6] and making a determination of the radial distribution of NII ions [7]. DIRBE has also contributed to interplanetary astronomy by providing an accurate phenomenological model of the Zodi Light from the interplanetary dust cloud [8].

Figure 1 shows the DIRBE annual average 100 micrometre map which is an example of the most detailed map data with highest contrast and largest dynamic range.

## 3. Ground Segment Architecture

The ground segment computer architecture consists of a VAXcluster linked by an Ethernet network bridged by a hardware-based repeater. It supports approximately 100 users in the daytime, production work at all hours, and system management and monitoring activities [9]. The HSC's serve 100 Gbytes of magnetic disk to the cluster, which consists of four mainframes and thirteen workstations. Interactive development and analysis work is done on the workstations which provide almost all the CPU power in the cluster. The mainframes are reserved for disk serving and batch processing. With the advent of truly high performance workstations, the I/O demands are also increased and disk serving has become a critical load to all but the most powerful mainframes. Two DECStation 5000 workstations are currently available and are linked to the VAXcluster using NFS. The data sets generated by the project pipelines are available to remote users and PCs through a data server and can be manipulated using IDL which is in widespread use on the VAX/VMS platforms.

## 4. Project Data Sets

The COBE satellite carries three experiments designed to make high precision measurements of the diffuse celestial background.

The detectors are stable and data sampling highly redundant. The observed sky is faint, low-contrast and smoothly variable except for one instrument (DIRBE) which sees stars at fixed map coordinates. FIRAS and DIRBE report glitches, many of which arise during passages over the South Atlantic Anomaly region.

The processed data currently totals 380 GB with an effective expansion factor of (4-16) over the raw data which depends on the instrument system. The project standard data sets number about 1000 and may be classed as sky maps, time-ordered data and time-tagged data. These data sets combine scientific with engineering data.

The Project Data Sets are required to represent data free of instrumental signature and the Analyzed Science Data Sets are intermediate to further scientific interpretation. The Astronomical Databases [10] contain external survey data converted to the COBE sky cube pixelization scheme [11],[12] at the resolution and beam pattern of the COBE instruments. The sky cube is an approximate equal-area projection on the sky of the faces of an inscribed cube. The equal-area property is ensured by the curvilinear coordinate system ruled on each cube face.

COBE data sets are directories of files. Intensity, spectral and polarimetric data are stored in area quadtree maps together with ancillary information. Offsets into each map corresponding to each level of resolution are stored in "index" files. Each pixel may contain one or more records with the same field structure. Data destined for the DIRBE experiment are stored at sky cube pixel level with 9 or more levels of resolution available in an image pyramid obtained by spatial averaging; data intended for FIRAS and DMR are stored at 6 or more levels of resolution. Data records are fixed length, defined by a Record Definition Language (RDL) file interface to the VAX Common Data Dictionary. RDL and its Record Definition Compiler were developed by the COBE project [9].

Figure 2 shows an example RDL for the DIRBE Daily File.

## 5. Data Compression Requirements

Data Compression is intended to simplify the task of systems management, data migration and recovery from catastrophic disk failures, reduce expenditure on storage devices and improve the data retrieval rate by a substantial factor dependent on the non-linear response of the saturated network.

The COBE Ground Segment Software System [9] consists of approximately 500 packages known as *facilities* which process the data in *pipelines* for each subsystem from raw telemetry to Project Data Sets. Access to the data is provided by the Data Management subsystem heavily dependent on a project-specific access system known as *COBEtrieve*.

Interviews with the Principal Investigators and Contract Leaders defined requirements as follows:

Provide compression transparently without changing the application software.

Compress instrument pipeline and science analysis data products to better than (16 to 50)%.

Process compressed data at a throughput not less than 90% of uncompressed data processing (possibly faster).

Preserve required accuracy of instrument housekeeping and scientific data (as judged by validators).

Exceed bitwise reliability of  $10^{-13}$  on average (flawless compression of 300 GB). Several times this factor is desirable.

Support full random access to file records.

Provide a capability to select specific classes of data for compression.

Preserve overlaps in separately-processed data segments.

Store search keys (time code, pixel address) in clear codes.

Provide a capability to select a compression scheme for each field of a data record.

Optimize choice of compression scheme combining *a-priori* with adaptive knowledge of data.

## 6. Implementation

Initial tests with public domain software (Unix-compress) and commercial PC-based hardware (Stacker, a product of Stac Electronics) demonstrated poor performance. The software was far too slow to keep up with the processing and Stacker compressed the DIRBE Daily Files (the largest archived files with the greatest retention time) by < 2%. Although the offlining of disk volumes provided by the FlashDAT 4mm tape device (a product of Winchester Technologies, Inc.) has been highly effective (factor of 4 improvement in data migration rate with a compression factor of 2), the requirements listed above necessitate customized software.

The following decisions were taken :

Create standalone, callable and embedded software interfaces.

Use existing fixed-length file record structure.

Use existing search algorithms to retrieve data.

Store compression parameters in file header without increasing the number of open files. This averts a resource lock-management problem in an already full file system.

Adopt an incremental build strategy: simple, well-trusted algorithms followed by powerful sophisticated methods.

Assess all algorithms on samples of all types of project data : [Quadtree Sky Map, Time-Ordered, Time-Tagged].

Optimize tradeoff between throughput and compression factor via overall measured storage savings.

Store data for medium-term via deeply-compressive but slow methods assisted by accelerator board (single files recovered in << 8 hours).

Offline project data via hardware-based compression methods.

Data shall not be delivered in compressed form to external users.

## **7. Compression System Design**

Since the data access is heavily dependent on *COBE*trieve and all the I/O system calls are localized, a natural solution is to embed compression software between the data management and I/O layers. This software compresses and decompresses data from the stored format to the fixed-length record structure understood by the data management software.

The writing of compressed data may be toggled via a system-wide logical name. The reading of compressed data is always enabled. The compression method specification is via command-line qualifiers which may be stored in the compressed file header and parsed to control the decompression of archive files. These qualifiers drive the command line, callable and embedded interfaces uniformly.

Currently, the record length and connect-time attributes of recognized standard data sets are stored in a VAX Datatrieve data base (DAFS). When a file is opened, this data base is queried and if the data set name and record length are matched, the data are accessed. Separately-processed time-overlapping data segments are stored in separate files but the data streams are merged based on the most recently-processed data from each segment.

Similarly, we may define a compression data base (CMPR) that specifies the command-line qualifiers (including the record length) which will be parsed to control the (de)compression of archive files. Since multiple compression method types and offset endpoints are defined for multiple offset ranges, this data may require updating on every change of data set Record Definition Language specification. Ideally, this information would have been provided by the scientist when the data sets were being designed.

The compression system permits full upwards and downwards compatibility with existing files and catalogs. If the record-length matches the entry in the DAFS data base, file is assumed

uncompressed. If the entry does not match, the file header is parsed for the decompression parameters. The compression parameters for standard data sets stored in the CMPR data base may not be overridden by users (the command-line qualifiers will be ignored) so the compression technique for a standard data set is under configuration control.

If a file is compressed from the DCL command level, a system-unique temporary file is used to store the data. If the compression is successful (a shorter file is created), the original file is replaced by the temporary file. All permanent attributes except the record-length are retained. Since the file name, version, extension, creation and modification dates are unchanged, the archive catalog need not be updated. Since the modification date is unchanged, VMS BACKUP software will not restore any offlined version of the same file, reducing the offline storage volume.

### **8. Compression Method Specification**

The compression methods so far envisaged consider a packet of successive records ("chunk") as an image to be compressed. The methods may require parameters (such as a range), positional information (matrix partition) and a specification of the number of records in the buffer. Although non-optimal, each block, delimited by a specified field offset range, is constrained to a fixed number of records in the buffer. The block may be scanned in column order ("transposed"), row or rectangular image and variable-length output is reformatted and re-aligned to fixed-length records. An optional list of reference filenames may be provided and a list of floating point parameters may be required.

The following generic compression schemes are provided :

- Field : data fields are compressed by re-quantization.
- Scanline : data in a "horizontal" or "vertical" range of scanlines is compressed by methods which consider correlations between adjacent data elements. The FULL vertical (time-series) scanline is compressed.
- Block : data in a non-overlapping, multiple range of offsets is compressed by methods which consider the correlations between neighboring elements. The operators may be causal, acausal or semi-causal in scanline order.

### **9. Compressed Data Record Structure**

The existing data management system is based on a fixed-length record structure with field offsets defined in an RDL file. The record-length and connect-time file attributes are stored in a database under Configuration Management control. Since many data compression methods generate variable-length records, it was necessary to devise a scheme permitting full random access without wasting storage on record filler bytes.

Since the time code and pixel address label fields are strictly monotonically increasing (except for certain data sets not destined for compression) this may be achieved as follows:

Figure 3 demonstrates the separate compression of field offset ranges for "packets" of fixed-length records with fixed-length output records supporting full random access by time code and pixel address. The status byte indicates whether the record is compressed or not and the "lookback" word points to the beginning of the output record. The shaded areas denote successive samples of data in pre-defined offset ranges. Subrecords are broken across the record boundary with the label fields deferred to the beginning of the next output record. In this manner, if the search finds a label value the "lookback" field refers to the start of the compressed data associated with that label. The result is that these fields are never split across record boundaries and no space is wasted.

A restriction is placed on the length of an output record that it must not exceed the length of the "lookback" field plus the length of the status field. The output record length is constrained to always exceed this value so no input record may span more than two output records. Any output record that exceeds this limit is transmitted in clear codes. Any compressed file larger than the original is transmitted in uncompressed form.

#### **10. Random Access to Compressed Data.**

The efficient search for matching time codes in large time-ordered files requires the insertion of an internal time code index list at predefined records in the uncompressed file. When a file is opened the first index list is read into virtual memory. If the desired record is not in the decompressed buffer, the bounding time codes are searched for in the index list to minimize the I/O. If the time codes are not found in the current index list, the next list is read into memory. The search uses a hunt and locate method, where the initial record is predicted from the average compression factor for the file, determined from the compressed file size and the number of uncompressed records multiplied by the uncompressed record length stored in the archive catalog. The exponential search is carried out until the time codes are bracketed when a binary search is used to locate the exact compressed record. The compressed record buffer is searched linearly for the matching time codes. The reduction in I/O by using index lists leads to an order of magnitude improvement in search time.

The search for matching pixel addresses in a sky map proceeds similarly except that an index file pointing to the first logical (uncompressed) record under a pixel is already available. Two lists of corresponding logical and physical (compressed) record numbers are stored in the file. In both cases, the index lists are highly compressible.

## 11. Compression Algorithms

The initial algorithmic toolbox will contain range quantization, run-length coding and zero suppression methods. The range quantization is an approximate method currently used by DIRBE which recognizes the sentinel values flagging noisy data.

Planned subsequent development includes nested Chebyshev polynomials (smoothly-variable data), a modified Huffman code, the Haar Transform followed by quadtree bit plane encoding, variants of the Lempel-Ziv-Welch substitution schemes with static codebooks, stochastic models such as the Autoregressive Integrated Moving Average (ARIMA) schemes and tree-structured Vector Quantization based on static codebooks. Since the data distribution is almost stationary with time, a static codebook may be stored on in memory for codebook-based algorithms. Usage of the Vector Quantization algorithms will depend on available resources and will probably be restricted to static archives.

## 12. Worked Example

The DIRBE experiment was operated with cryogenic cooling for 41 weeks, creating 80MB per day for a total of 5.5GB processed data.

Clearly, this RDL was devised with each field carefully specified for scientific usage and it is not necessary to make minimalist assumptions about the nature of the data. This RDL specifies a mixture of scientific and engineering data and some fields must be transmitted in clear codes (search labels), exactly (flags), approximately (photometry) or are noisy and hence incompressible (e.g. pixel subposition).

The records are 140 bytes long and in quadsphere sky map format [10]. The label field is the "Pixel\_No" which is referenced explicitly in the user software as a pixel-number-offset argument to the access software. There are 16 floating-point photometric bands.

Direct usage of "Unix-compress" leads to a compression factor of 25% which takes several hours to compress one sky map on a workstation.

DIRBE has already decided that a logarithmic range compression scheme which sentinelizes glitchy data (flagged in a previous pipeline process) is sufficient to convert the floating-point photometry to 16 bit integers on a field-by-field basis. Further compression may be achieved particularly for data which are not glitchy (Glitch\_Flags) or taken in a particle radiation zone (Radiation\_Cont). This represents about 75% of the data. This compressible data may be vector quantized with a suitable codebook derived (perhaps) by the Linde-Buzo-Gray algorithm based on a training set extracted from a typical daily file. A normalized codebook would be the most flexible. At best, this approach would yield ~ 2 bytes per array of highly-correlated photometry bands.

The ratio of the daily photometry to the annual average value under a pixel is expected to have reduced dynamic range and be even more compressible. The "pixel\_no" fields and the ancillary angles between the DIRBE boresight and celestial objects which vary slowly under a pixel are ~50% compressible via a Modified Huffman Code in vertical scan mode. The overall compression factor is about 50%.

In another example, one FIRAS facility accesses time-ordered data via extensive keyed-read operations which involve searches which currently create the largest single network load. The files are approximately 16MB consisting of 8 byte time codes together with ~ 10 000 bytes of data. Each search step (there are typically about 6 per keyed read) reads the whole record to locate the time code. The compressed data which contains the internal time code list may be searched ~ 30 times faster as the total I/O is reduced to 10% of its original value.

### **13. Validation and Testing**

Software quality has been assured by regression testing in an independent environment to ensure that goals of functionality, accuracy and performance have been met. Code inspection has been used to ensure the robustness and maintainability of the code and documentation.

The in-house validation team will provide quality assurance for the compressed data products using the same formal project accuracy requirements as for original data.

Tests of file migration to/from all available media (including 4mm and 8 mm magnetic tape, magnetic and optical disks and 9-track tape) indicate that the compressed data files are fully compatible with VMS BACKUP and COPY software and that the project-specific data migration software facility is effective with compressed data.

### **14. Summary and Conclusions**

A general approach incorporating scientific knowledge seems appropriate for the Space and Earth Science Data Compression application. Inline data compression techniques developed for the COBE project may help the project achieve its goals and be useful to other workers in this growing field.

### **15. Recommendations for Future Development**

Compression functions should be specified at the same time the data sets are defined. An optimal implementation may consider the data as a linked list of object classes for each data field which specify overloaded (de)compression functions invoked in the constructor for each class.

A Data Compression Designer Expert System could capture the knowledge of domain experts and recommend appropriate functions.

## Acknowledgements

The COBE data analysis is managed by the Goddard Space Flight Center for NASA's Office of Space Science and Applications.

COBE is a team effort and more than 1000 individuals have contributed to its success.

## References

- [1] J. A. Storer, M. Cohn (eds.), Proc. Data Compression Conf. 1992, held at Snowbird, Utah, April 24-27, 1992.
- [2] N. W. Boggess et al., "The COBE Mission; Its Design and Performance Two Years after Launch", Ap. J. 1 October, 1992.
- [3] J. C. Mather et al., "The COBE FIRAS Measurement of the Cosmic Microwave Background Spectrum", Ap. J. , 1993 (in preparation).
- [4] M. G. Hauser et al., "The Diffuse Infrared Background : COBE and other results", in "After the First Three Minutes", AIP Conf. Proc. 222, 1991.
- [5] J. Weiland et al., "COBE/Diffuse Infrared Background Experiment (DIRBE) Observations of the Milky Way Galactic Bulge", Ap. J. Lett., 1993 (in preparation).
- [6] T. Sodroski et al., "Large Scale Physical Conditions in the Interstellar Medium from DIRBE Observations", Ap. J. Lett., 1993 (in preparation).
- [7] E. L. Wright et al., "Preliminary Spectral Observations of the Galaxy with a 7 degree beam by the Cosmic Background Explorer (COBE)", Ap. J., 21 October, 1991.
- [8] W. Spiesmann et al., "Near and Far Infrared Observations of Interplanetary Dust Bands from COBE", Ap. J. Lett., 1993 (in preparation).
- [9] E. S. Cheng, "COBE: The Software", Proc. First Astronomical Data Analysis Software and Systems Conf., held at Tucson, Arizona, Nov 8-11, 1993.
- [10] I. Freedman, A. C. Raugh, E. S. Cheng, "The COBE Astronomical Databases", Proc. First Astronomical Data Analysis Software and Systems Conf., op. cit.
- [11] R. A. White, S. W. Stemwedel, "The Quadrilateralized Spherical Cube and Quad-Tree for All Sky Data", Proc. First Astronomical Data Analysis Software and Systems Conf., op.cit.
- [12] I. M. O'Neill, R. E. Laubscher, "Extended Studies of a Quadrilateralized Spherical Cube Earth Data Base", NEPRF Tech. Rep. 3-76 (CSC).

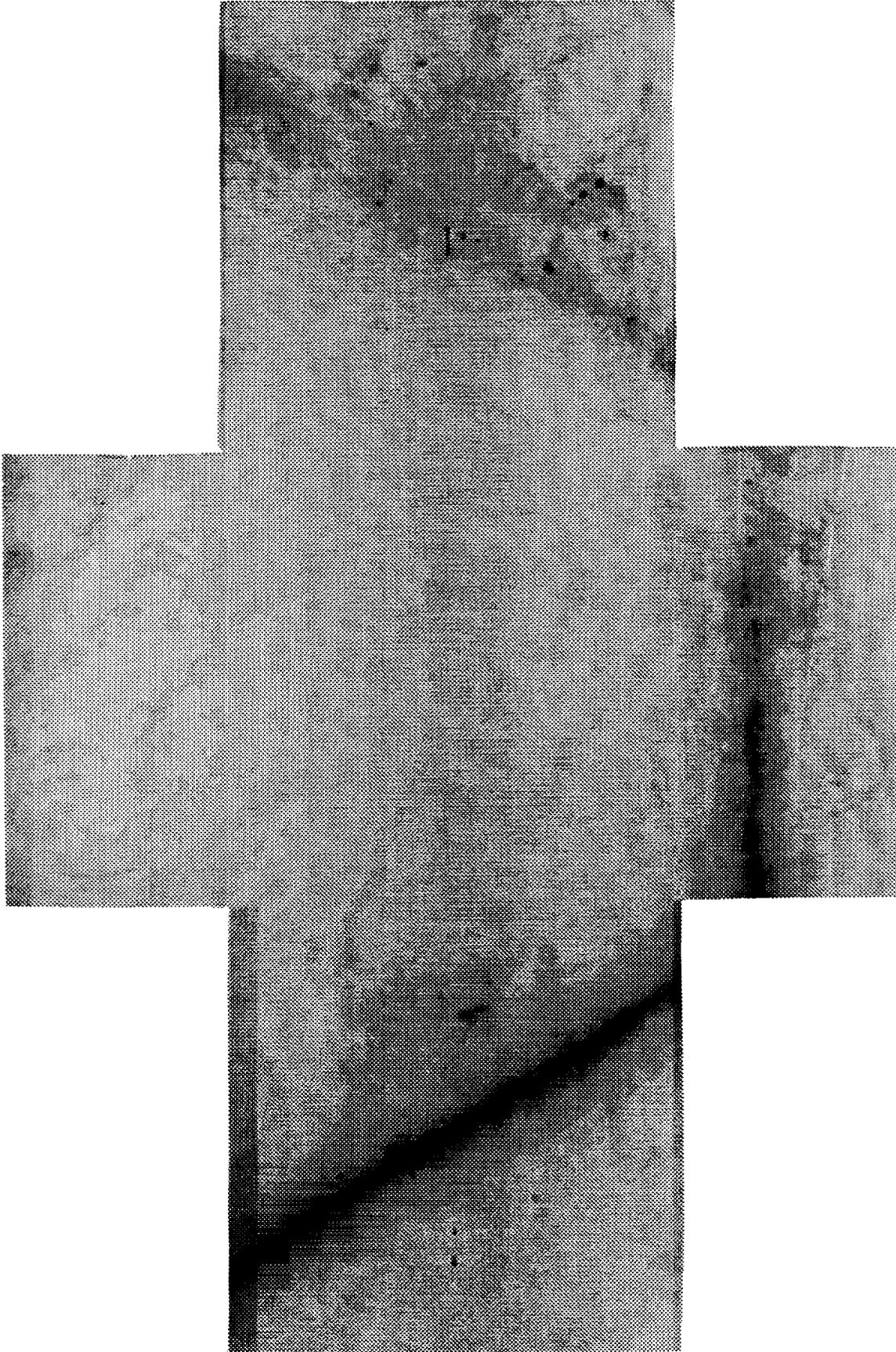


Figure 1. Annual Average 100  $\mu\text{m}$  map projected on quadisphere.

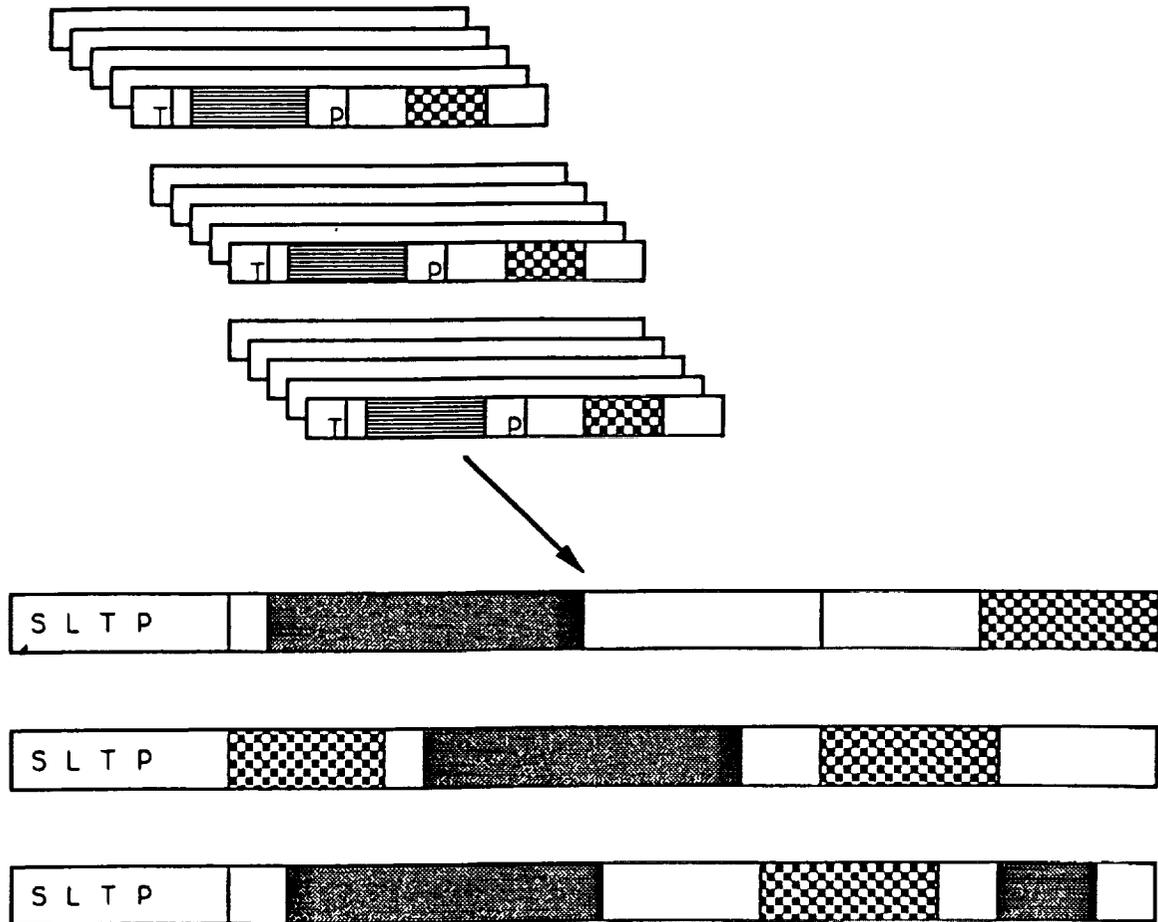
RECORD BCI\_CIRSSM BCI\_CIRSSM ! Complete IRS Sky Maps

Offset	Length	Description
0	8	SCALAR Time /ADT!Time of middle of observation.
8	4	SCALAR Pixel_no /LONG!Pixel number of observation
12	64	ARRAY Photometry /FLOAT/DIM=16!Detector observations
76	1	SCALAR Approach_vector /BYTEU!Forward looking = 1 !Backward looking = 2 (referenced to SC velocity)
77	1	SCALAR Pixel_subpos /BYTEU!Sub-pixel containing DIRBE LOS
78	4	SCALAR Next_obs /LONG!Pixel number of next observation
82	4	SCALAR Prev_obs /LONG! previous
86	4	ARRAY Sun_re_BS /WORD/DIM=2 !Word 1: Solar elongation. !Word 2: Relative azimuth of sun
90	2	SCALAR SC_Axis_re_Zenith /WORD!Angle between COBE -X axis !and the zenith vector
92	2	SCALAR BS_re_Zenith /WORD!Angle between DIRBE boresight and
94	2	SCALAR BS_re_Horiz /WORD!Angle between earth horizon and ! DIRBE boresight.
96	2	SCALAR SC_Axis_re_Vel /WORD!Angle of COBE -X axis relative !to velocity vector
98	2	SCALAR BS_re_Vel /WORD !Angle between DIRBE boresight !and S/C velocity vector
100	2	SCALAR Azimuth_re_Vel /WORD
102	6	ARRAY Attack_vector /WORD/DIM=3
108	2	SCALAR FOV_Azimuth /WORD
110	2	SCALAR Longitude /WORD
112	2	SCALAR Latitude /WORD
114	2	SCALAR Altitude /WORD
116	3	ARRAY Mag_Field /BYTE/DIM=3
119	4	ARRAY Moon_re_BS /WORD/DIM=2
123	2	SCALAR Sun_Moon_Angle /WORD
125	2	SCALAR Moon_Distance /WORD
127	4	ARRAY Jupiter_re_BS /WORD/DIM=2
131	4	ARRAY Earth_light_cont /WORD/DIM=2
135	1	SCALAR Pixel_subsubpos /BYTEU
136	1	SCALAR Radiation_cont /BYTEU
137	2	SCALAR Glitch_Flags /WORDU
139	1	SCALAR ATT_Flags /BYTEU
140		END_RECORD

TOTAL LENGTH OF RECORD: 140 BYTES  
 TOTAL NUMBER OF FIELDS: 29

Figure 2. Record Definition Language for DIRBE Daily File.

# FRAMEWORK



Separate compression of field offset ranges for "packets" of fixed-length records with fixed-length output records supporting full random access by time-tag and pixel number. The status byte indicates whether the record is compressed or not and the "lookback" word points to the beginning of the output record. The shaded areas denote successive samples of data in pre-defined offset ranges. The notation "S L T P" denotes status, lookback, time-tag and pixel number.

**Figure 3. Separate compression of field offset ranges.**



**PROPOSED DATA COMPRESSION SCHEMES FOR  
THE GALILEO S-BAND CONTINGENCY MISSION \***

Kar-Ming Cheung      Kevin Tong  
Communications Systems Research  
238-420  
Jet Propulsion Laboratory  
4800 Oak Grove Drive  
Pasadena, CA 91109

**Abstract.** The Galileo spacecraft is currently on its way to Jupiter and its moons. In April 1991, the high gain antenna (HGA) failed to deploy as commanded. In case the current efforts to deploy the HGA fails, communications during the Jupiter encounters will be through one of two low gain antenna (LGA) on an S-band (2.3 Ghz) carrier. A lot of effort has been and will be conducted to attempt to open the HGA. Also various options for improving Galileo's telemetry downlink performance are being evaluated in the event that the HGA will not open at Jupiter arrival. Among all viable options the most promising and powerful one is to perform image and non-image data compression in software onboard the spacecraft. This involves in-flight re-programming of the existing flight software of Galileo's Command and Data Subsystem processors and Attitude and Articulation Control System (AACS) processor, which have very limited computational and memory resources. In this article we describe the proposed data compression algorithms and give their respective compression performance.

The planned image compression algorithm is a  $4 \times 4$  or an  $8 \times 8$  multiplication-free integer cosine transform (ICT) scheme, which can be viewed as an integer approximation of the popular discrete cosine transform (DCT) scheme. The implementation complexity of the ICT schemes is much lower than the DCT-based schemes, yet the performances of the two algorithms are indistinguishable.

The proposed non-image compression algorithm is a Lempel-Ziv-Welch (LZW) variant, which is a lossless universal compression algorithm based on a dynamic dictionary lookup table. We developed a simple and efficient hashing function to perform the string search.

## 1. Introduction

The Galileo spacecraft, which was launched in Oct 1989, is now on its way to Jupiter. Its mission includes releasing a probe into the Jovian atmosphere, Io flyby, probe data capture and relay, Jupiter orbital insertion, and 10 satellite encounters (Ganymede, Callisto, Europa). The Galileo project involves over 20 years of effort. In April 1991, when the spacecraft first flew by Earth, the Galileo team commanded the spacecraft to open the 1.8m high-gain antenna (HGA). However, the HGA failed to completely deploy. All indications are that 3 of the 18 ribs are stuck to the antenna's central tower. Several unsuccessful attempts have been made to free the stuck ribs. A major effort is planned for December 1992 to perform hammering or pulsing of the deployment motor to try to free the ribs. If the HGA fails to deploy, the only way to communicate between

---

\* The research described in this paper was carried out by Jet Propulsion Laboratory, California Institute of Technology, under a contract with National Aeronautics and Space Administration

Earth and the spacecraft is through the use of one of the two low gain antennas. If the current configuration (ground and spacecraft) remains unchanged, the telemetry data rate will be 10 bits per second at Jupiter arrival (1995), compared to the expected data rate of 134 kbits per second in the HGA configuration. The amount of data that can be returned would be drastically reduced.

A study [1] was conducted from December 1991 through March 1992 to evaluate various options for improving Galileo's telemetry downlink performance in the event that the HGA does not open by Jupiter arrival. Among all viable options the most promising and powerful one is to perform image and non-image data compression in software onboard the spacecraft. This involves in-flight re-programming of the existing flight software of Galileo's Command and Data Subsystem (CDS) processors and the Attitude and Articulation Control System (AACS) processor, which has severely limited computational and memory resources. The software has to be compact and computationally simple. A lossy image compression scheme is proposed that can give a wide range of rate-distortion trade-off for the image data, represents over 70% of the data to be returned by the mission. The rest of the data comes from various spacecraft instruments. This can either be compressed by using instrument-specific compression schemes or by using a proposed lossless universal compression algorithm. In this article we describe the proposed image compression scheme and the universal lossless compression algorithm and give their respective compression performances.

The proposed image compression algorithm is a  $4 \times 4$  or an  $8 \times 8$  multiplication-free integer cosine transform (ICT) [2], which was first proposed by Cham. The ICT can be viewed as an integer approximation of the popular discrete cosine transform (DCT) scheme. The  $8 \times 8$  multiplication-free ICT will be implemented in software using the more powerful AACS processor and the  $4 \times 4$  ICT will be implemented in software using several CDS processors as backup. The ICT schemes have much lower implementation complexity and give indistinguishable performances when compared to the DCT schemes.

The proposed non-image compression algorithm is a Lempel-Ziv-Welch (LZW) variant [3], which is a lossless universal compression scheme. Due to the severe limitations of the CDS processors, we cannot use the more sophisticated existing hashing functions [3]. We developed a simple and efficient hashing algorithm to perform string search. This hashing function uses a total 1802 bytes of memory for a codebook of size 512 bytes, and requires on the average only 3 16-bit comparisons per input byte.

The communication scenario described in this article is unique. Rather than a typical data compression paradigm as in industry where a sophisticated encoder and simple decoder are desired, the Galileo HGA anomaly situation requires a very simple compressor onboard the spacecraft. The decompressor, which is on the ground, can be reasonably complex. Many of the compression techniques described in this article are not novel and are modifications and enhancements of some existing algorithms to adapt to the HGA anomaly scenario. The main goal is to simplify the onboard compressor implementation.

The rest of this article is organized as follows: Sections I, II, III, IV, V, VI and VIII describe the ICT lossy image compression scheme. A more in-depth discussion of the relationship between ICT and DCT is given in Section II. The interplay between the orthonormal transform stage (any orthonormal transform, not just DCT) and the quantization stage is explored in Section III. The mathematical properties and a general construction scheme for the multiplication-free ICT matrix are given in Section IV. A general construction procedure of ICT matrix is described in SECTION

V. Examples of  $4 \times 4$ ,  $8 \times 8$ , and  $16 \times 16$  ICT matrices are given in Section VI. The rate-distortion performance of the ICT schemes for the Galileo S-Band Contingency Mission ( $4 \times 4$  and  $8 \times 8$ ) is described in Section VII. Section VIII gives an overview of the LZW algorithm. Section IX describes the LZW scheme we used in the Galileo LGA mission. Section X describes the novel features of the Galileo LZW implementation.

## 2. DCT Versus ICT

The discrete cosine transform (DCT) is regarded as one of the best transform techniques in image coding. Its independence from the source data and the availability of fast transform algorithms make the DCT an attractive candidate for many practical image processing applications. In fact the ISO/CCITT standards of image processing in both still-image and video transmissions includes the two-dimensional DCT as a standard processing component in many applications. For still-image compression, the transform-based scheme consists of three stages: the data transform stage, the quantization stage and the entropy-coding stage. For video compression, an additional motion-compensation stage with feed-back is included. The enormous popularity of the DCT in image compression provides the driving force for researchers to develop efficient hardware and software implementations for the DCT.

The commercial acceptance of the emerging JPEG and MPGE Standards, which uses an  $8 \times 8$  block DCT has created a need for an efficient DCT algorithm. A lot of effort has been devoted to the pursuit of reducing the computational complexity of the DCT. New algorithms have already been proposed[6][7]. The idea of incorporating the scale-factors of the transform process as part of the scalar quantizer has also been suggested in the recent literature [6][7]. All these efforts gear towards reducing the total number of floating-point or fixed-point multiplications and additions used, with the emphasis on reducing the number of multiplications.

Recently Cham [2] took a different approach and proposed a new 8-point transform called the integer cosine transform (ICT). ICT requires only integer multiplications and additions, making it much simpler to implement than the DCT. An ICT chip was fabricated and was proven to be efficient in both silicon area and speed. The elements in an ICT matrix are all integers, with sign and magnitude patterns that resemble those of the DCT matrix. The similarity of the ICT matrix to the DCT matrix, together with the orthogonality property of the ICT ( $CC^t = \Delta$ , where  $C$  is an ICT matrix and  $\Delta$  is a diagonal matrix), guarantee that the ICT as well as its inverse possess the same transform structure as the DCT, thus allowing the use of some fast DCT algorithms to compute a fast ICT [2].

Although the 8-point ICT proposed by Cham performs remarkably well, it is quite ad hoc. In this article we put ICT into a more formal mathematical setting and generalize Cham's idea to any  $N$ -point ICT. The mathematical properties of orthonormal transforms including ICT are investigated in the following sections. Since ICT is separable and the extension of the one dimensional ICT to two dimensions is straight-forward, this article focuses on the one dimensional case.

### 3. Orthonormal Versus Orthogonal Separable Transforms

An  $N \times N$  1-D matrix  $M$  is said to be orthonormal if and only if  $MM^T = I$ , where  $I$  is the identity matrix. An  $N \times N$  1-D matrix  $C$  is said to be orthogonal if  $CC^T = \Delta$ , where  $\Delta$  is a diagonal matrix. It can be shown from basic linear algebra that for any  $N \times N$  orthogonal matrix  $C$ , there exists an  $N \times N$  orthonormal matrix  $M$  and an  $N \times N$  diagonal matrix  $\Delta$  such that  $M = \sqrt{\Delta}C$ . It can further be shown that  $C^{-1} = C^T\Delta$  and  $M^{-1} = C^{-1}\sqrt{\Delta}^{-1} = C^T\Delta\sqrt{\Delta} = C^T\sqrt{\Delta}$ .

The corresponding 2-D  $N^2 \times N^2$  orthonormal separable transform matrix is

$$M \otimes M = (\sqrt{\Delta}C \otimes \sqrt{\Delta}C) = (\sqrt{\Delta} \otimes \sqrt{\Delta})(C \otimes C), \quad (1)$$

where  $X \otimes Y$  denotes the tensor product of the matrix  $X$  with  $Y$ , and the corresponding 2-D  $N^2 \times N^2$  orthonormal inverse transform matrix is

$$(M \otimes M)^{-1} = (M^{-1} \otimes M^{-1}) = (C^T\sqrt{\Delta} \otimes C^T\sqrt{\Delta}) = (C^T \otimes C^T)(\sqrt{\Delta} \otimes \sqrt{\Delta}), \quad (2)$$

The matrix  $\sqrt{\Delta} \otimes \sqrt{\Delta}$  is diagonal. Therefore when the 2-D orthonormal transform  $M \otimes M$  is followed by quantization, the diagonal matrix  $\sqrt{\Delta} \otimes \sqrt{\Delta}$  can be absorbed in the quantization stage and, only the product by the orthogonal matrix  $C \otimes C$  is computed in the transform stage. Similarly on the decoder side,  $\sqrt{\Delta} \otimes \sqrt{\Delta}$  can be absorbed in the de-quantization stage, and the  $N^2$  output samples from the de-quantizer are multiplied by the orthogonal matrix  $C^T \otimes C^T$ . The fusion of the scaling factors of the transform (inverse) transform stage into the quantization (de-quantization) stage does not require additional computation, since division operations have to be performed in the quantization process anyway. An example of a quantization stepsize template that corresponds to the all-one uniform quantization template for an  $8 \times 8$  ICT is given in Figure 2. A more detailed discussion on incorporating the scale-factors of the transform process as part of the scalar quantizer can be found in [7]. This relaxation of the orthonormal requirement to orthogonal requirement play a crucial role in allowing one to "integerize" a transform coding scheme as we will see in the next section.

### 4. Mathematical Properties of ICT

ICT and DCT are closely related. Let  $C$  and  $A$  be the respective ICT and DCT  $N \times N$  matrices.  $A = [a_{kn}]$  is an orthonormal matrix (i.e.  $AA^t = I$ ) defined as follows:

$$\begin{aligned} a_{kn} &= \frac{1}{\sqrt{N}} & k = 0, 0 \leq n \leq N-1 \\ &= \sqrt{\frac{2}{N}} \cos \frac{\pi(2n+1)k}{2N} & 1 \leq k \leq N-1, 0 \leq n \leq N-1 \end{aligned} \quad (3)$$

Using  $A$  as a template, the ICT matrix  $C = [c_{kn}]$  is an orthogonal matrix (i.e.  $CC^t = \Delta$ , where  $\Delta$  is a diagonal matrix) with the following properties:

1. Integer property:  $c_{kn}$  are integers for  $0 \leq k, n \leq N-1$ .
2. Orthogonality property: Rows (or columns) of  $C$  are orthogonal.

3. Relationship with DCT: (i)  $\text{sgn}(c_{kn}) = \text{sgn}(a_{kn})$  for  $0 \leq k, n \leq N - 1$ . (ii) If  $a_{kn} = a_{st}$ , then  $c_{kn} = c_{st}$  for  $0 \leq k, n, s, t \leq N - 1$ .

The integer property eliminates real multiplication and real addition operations. The orthogonality property assures that the inverse ICT has the same transform structure as the ICT. Notice that  $C$  is only required to be orthogonal, but not orthonormal. However, any orthogonal matrix can be made orthonormal by multiplying it by an appropriate diagonal matrix. This operation can be incorporated in the quantization (dequantization) stage of the compression (decompression) scheme, thus sparing the ICT (inverse ICT) transform from floating-point operations, and at the same time preserving the same transform structure as in the floating-point DCT (inverse DCT). The relationship between ICT and DCT guarantees efficient energy packing and allows the use of some fast DCT technique for the ICT.

## 5. A General Procedure to Construct an ICT Matrix

A general procedure to construct an  $N \times N$  ICT matrix is presented in this section. For any  $N \times N$  ICT matrix, this construction is done on the ground prior to implementation. The DCT matrix is used as a template to generate an ICT matrix. The procedure is described as follows:

1. Generate the  $N \times N$  DCT matrix  $A$ .
2. Construct an  $N \times N$  matrix  $C$  by substituting the  $N$  possible absolute values in  $A$  with  $N$  symbols, and preserve the signs of the elements in  $A$ .
3. Evaluate  $CC^t$ , and generate a set of independent algebraic equations which forces  $CC^t$  to be a diagonal matrix.
4. Find a set of  $N$  numbers which satisfies the set of algebraic equations generated in part 3.

Since for a given  $N$ , there are  $N(N - 1)$  non-diagonal elements in  $C$ , part (3) gives  $N(N - 1)/2$  quadratic equations. This set of equations is too large to be handled easily except for small  $N$ . The most tedious part of the above procedure is part 4, that is finding  $N$  integers satisfying the set of non-linear algebraic equations generated in part 3. By using advanced symbolic manipulation tools like *Mathematica* [8], the effort to generate the set of algebraic equations in part 3 and solving them in part 4 can be greatly reduced. In fact *Mathematica* was used in an interactive manner to generate a  $4 \times 4$ , an  $8 \times 8$  and a  $16 \times 16$  ICT matrices as described in the next section.

In order to obtain good compression performance one requires the set of  $N - 1$  integers to have a similar magnitude profile to the  $N - 1$  floating-point elements of  $A$ . Furthermore, if the multiplication-free property is desired, one has to restrict the set of  $N$  integers to be small integers, so that any multiplications with the matrix elements can be replaced by a small number of adds and shifts.

## 6. Examples of ICT Matrix Construction

Using the construction techniques described in the previous section, we generated a  $4 \times 4$  (Figure 1), an  $8 \times 8$  (Figure 2), and a  $16 \times 16$  (Figure 3) ICT matrices. The  $4 \times 4$  ICT matrix has

elements which are powers of 2. The  $8 \times 8$  ICT matrix is the same as the example given in Cham's paper [2], whose elements are either powers of 2, or are sums of two powers of 2.

## 7. Compression Performance of the ICT Schemes

We applied our implementation of the  $4 \times 4$  and the  $8 \times 8$  ICT schemes for the Galileo S-Band Contingency Mission. We compressed a typical planetary image *miranda* (moon of Uranus). For the purpose of comparison, we also compressed the same image using the JPEG schemes. The root-mean-square-error (RMSE) versus compression ratio performances of these schemes on *miranda* are given in Figure 4. These simulation results indicate that the difference in rate-distortion performance resulting from using the floating-point DCT or the ICT is unnoticeable.

The ICT schemes are also being considered for compression of non-image data like the multi-spectral plasma wave spectrometer (PWS) data. We compressed some typical PWS data files by a factor of 10, which results in lossy reconstructed images that can still be useful for PWS analysis.

## 8. LZW Overview

The universal lossless LZW algorithm used in this mission is based on the algorithm proposed by Terry A. Welch[3].

The LZW algorithm is an adaptive compression scheme which converts a variable length string into a fixed length string. The algorithm is adaptive in the sense that it uses a dynamic lookup dictionary table. The table is dynamic because it initially starts with an empty table of symbol strings and the algorithm fills this table during the compression and decompression process. The table is thus adapted to the incoming data. Because of this adaptation, the algorithm requires no prior information on the data characteristics of the incoming data.

The LZW implementation of the compression and decompression scheme is based on Welch's paper[3] with modifications to handle multiple dictionary tables, a more efficient search algorithm and the ability to detect certain errors. However, it must be noted that there is an error in the decompression algorithm described in Welch's paper. If followed exactly, the decompressed data will be garbled at random points. This error is located in the "special case" part of decompression algorithm defined in Welch's paper. Instead of a direct output of the decoded final character, this character should be pushed onto the stack.

Our contribution in this paper is to develop an LZW scheme that is efficient in terms of speed and compression performance and at the same time satisfies the stringent computation and memory constraints of the spacecraft.

## 9. LZW Algorithm

The LZW algorithm is organized around a translation table, referred to as a dictionary table that maps strings of input symbols into fixed length codes. In this particular mission, the code size used is 9-bits, which translates into a table size of 512. The dictionary is used as a lookup table in

both the compression and decompression and is generated as the data is being processed. If the required information is not in the present state of the table, a new entity is added to the table, thus a dynamic lookup table.

The compression speed is very sensitive to the search of the dictionary table in the main loop of the LZW algorithm. The search is used to determine if the required information is in the table. Since the entire LZW algorithm is based on the state of this table, it is important to develop a fast search routine that is also efficient in memory usage because of the memory constraint of the spacecraft.

### **9.1 The Dynamic Lookup Dictionary Table**

The size of the dictionary has a direct bearing on the memory requirements and execution time of the implemented program. The proposed dictionary size for this mission is 512. This number is a compromise between optimal compression and the memory constraint on board the spacecraft. The increase in dictionary size from 512 to 1024 does not produce a great enough compression gain to justify choosing the larger dictionary size.

## **10. Features of the Galileo LZW Implementation**

The LZW algorithm was implemented with features that were not discussed in Welch's paper. The implementation can concurrently compress multiple independent streams of data using multiple dictionaries while using the a minimal amount of memory without compromising execution time.

### **10.1 Multiple Dictionary Tables**

The multiple dictionary table feature was added because the spacecraft transmits different types of data requiring lossless compression. Examples of these types of data are telemetry, engineering and instrument data. Using the multiple dictionary approach, it is possible to segregate these data streams without requiring the compression algorithm to finish up on one stream and start another table. The program can switch back and forth between the data streams and use the dictionary table that is assigned to that data stream.

### **10.2 The Hashing Algorithm**

The search portion of the LZW algorithm is the most time consuming, thus it was necessary to design a search procedure that was both efficient in memory and in execution time. We employ a simple yet efficient hashing algorithm to perform the search. Normal implementation of hashing uses dynamic memory and a linked list, but in our implementation, two fixed arrays are used. This is to save memory and to save overhead time in keeping track of the linked list using dynamic memory. The size of the first array is equal to the dictionary size and the second array would equal the difference between the dictionary size and the alphabet size. Thus it would require one array of size 256 and the second array of size 512 for a dictionary size of 512. See Table 1 for hashing performance.

### 10.3 LZW Performance on Near-Infrared Mapping Spectrometer (NIMS) Data

We have obtained NIMS data produced by Galileo to test the performance of the LZW implementation. See Table 2.

### 10.4 LZW Performance on Selected Text Data

A text file was produced and used to show the performance of the LZW algorithm with various table sizes. (See Table 3.) From the performance, we can see that the table size proposed is a good compromise between optimal compression and memory usage.

### Acknowledgment

The authors would like to thank F. Pollara (Section 331) for his input to this work, and S. Dolinar (Section 331) and L. Swanson (Section 331) for their constructive comments. The authors would also like to thank T. Brady (Section 348) for his assessment on the implementation feasibility of the compression schemes described in this article. Particular thanks to the Deep Space Network Advance System Program whose support makes this project possible.

### References

- [1] L. Deutsch and J. Marr, "Galileo S-Band Mission Study, Final Report," JPL Publication, Mar 1992.
- [2] W. Cham, "Development of Integer Cosine Transforms by the Principle of Dyadic Symmetry," *IEE Proceedings*, Vol. 136, Pt. I, No. 4, August 1989.
- [3] T. Welch, "A Technique for High performance Data Compression," *Computer*, June 1984.
- [4] Joint Photographic Experts Group (JPEG) Draft Technical Specification (Rev. 5), ISO/CCITT, January 15, 1990.
- [5] Moving Picture Experts Group (MPEG) Draft Technical Specification, ISO/CCITT, May, 1991.
- [6] Y. Arai, T. Agui, M. Nakajima, "A Fast DCT-SQ Scheme for Images," *Trans. of the IEICE*, vol.E 71, pp.1095-1097, November, 1988.
- [7] E. Feig and S. Winograd, "Fast Algorithms for the Discrete Cosine Transform," submitted to the IEEE Transactions on Acoustics, Speech, and Signal Processing.
- [8] S. Wolfram, *Mathematica : A System for Doing Mathematics by Computer*, Addison-Wesley Publishing Company, 1988.

1	1	1	1
2	1	-1	-2
1	-1	-1	1
1	-2	2	-1

**Figure 1 a 4 x 4 ICT Matrix**

1	1	1	1	1	1	1	1
5	3	2	1	-1	-2	-3	-5
3	1	-1	-3	-3	-1	1	3
3	-1	-5	-2	2	5	1	-3
1	-1	-1	1	1	-1	-1	1
2	-5	1	3	-3	-1	5	-2
1	-3	3	-1	-1	3	-3	1
1	-2	3	-5	5	-3	2	-1

**Figure 2a an 8 x 8 ICT Matrix**

8	25	18	25	8	25	18	25
25	78	56	78	25	78	56	78
18	56	40	56	18	56	40	56
25	78	56	78	25	78	56	78
8	25	18	25	8	25	18	25
25	78	56	78	25	78	56	78
18	56	40	56	18	56	40	56
25	78	56	78	25	78	56	78

**Figure 2b the Quantization Template of 2a**

1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
42	38	37	32	22	19	10	4	-4	-10	-19	-22	-32	-37	-38	-42
10	9	6	2	-2	-6	-9	-10	-10	-9	-6	-2	2	6	9	10
38	22	4	-19	-37	-42	-32	-10	10	32	42	37	19	-4	-22	-38
2	5	-5	-2	-2	-5	5	2	2	5	-5	-2	-2	-5	5	2
37	4	-32	-38	-10	22	42	19	-19	-42	-22	10	38	32	-4	-37
9	-2	-10	-6	6	10	2	-9	-9	2	10	6	-6	-10	-2	9
32	-19	-38	4	42	10	-37	-22	22	37	-10	-42	-4	38	19	-32
1	-1	-1	1	1	-1	-1	1	1	-1	-1	1	1	-1	-1	1
22	-37	-10	42	-4	-38	19	32	-32	-19	38	4	-42	10	37	-22
6	-10	2	9	-9	-2	10	-6	-6	10	-2	-9	9	2	-10	6
19	-42	22	10	-38	32	4	-37	37	-4	-32	38	-10	-22	42	-19
5	-2	2	-5	-5	2	-2	5	5	-2	2	-5	-5	2	-2	5
10	-32	42	-37	19	4	-22	38	-38	22	-4	-19	37	-42	32	-10
2	-6	9	-10	10	-9	6	-2	-2	6	-9	10	-10	9	-6	2
4	-10	19	-22	32	-37	38	-42	42	-38	37	-32	22	-19	10	-4

**Figure 3 a 16 x 16 ICT Matrix**

Miranda

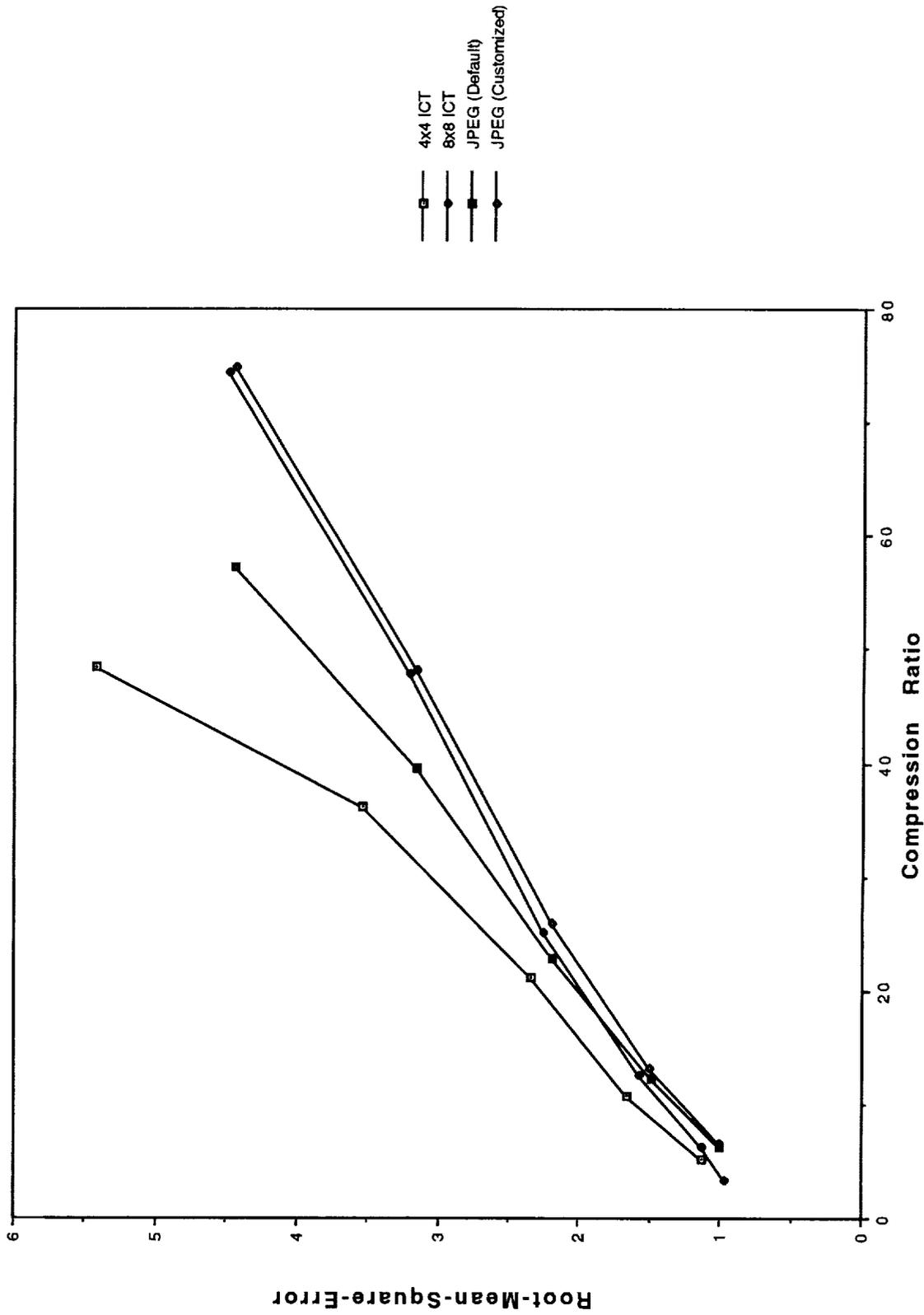


Figure 4 Compression Performances

**TABLE 1 Hashing Comparison**

Data files used are 256x256 = 65536 bytes planetary image files.

**Compares Per Input Byte (Table 1)**

File	Hash Algorithm Search		Sequential Search	
	512	1024	512	1024
d1	3.2436	3.7604	76.0222	184.5167
f2	2.5246	3.0067	73.0339	178.2090
h2	2.5816	3.1000	74.2097	181.3555
l2	3.6721	4.8428	93.6538	233.7537

**TABLE 2 Nims Data Performance**

**COMPRESSION RATIO (Table 2)**

Data Orientation	Table Size = 512	Table Size = 1024
Horizontal Scan	2.60	2.69
Vertical Scan	2.59	2.63
Mirror Scan	2.27	2.35
Original Data	2.45	2.51

**TABLE 3 Text Data Performance**

Sample Text File Size = 5390 bytes (Table 3)

Table Size	Compression Ratio	Tables Used
512	1.36	14
1024	1.52	4
2048	1.59	2



## DATA COMPRESSION FOR THE CASSINI RADIO AND PLASMA WAVE INSTRUMENT

W.M. Farrell, Code 695, NASA/Goddard Space Flight Center, Greenbelt, MD 20771, U.S.A.

D.A. Gurnett, D.L. Kirchner and W.S. Kurth, Department of Physics  
and Astronomy, The University of Iowa, Iowa City, IA 52242, U.S.A.

L.J.C. Woolliscroft, Department of Automatic Control and Systems  
Engineering, The University of Sheffield, Sheffield, S1 3JD, U.K.

**Abstract.** The Cassini Radio and Plasma Wave Science experiment will employ data compression to make effective use of the available data telemetry bandwidth. Some compression will be achieved by use of a lossless data compression chip and some by software in a dedicated 80C85 processor. A description of the instrument and data compression system are included in this report. Also, the selection of data compression systems and acceptability of data degradation is addressed.

### 1. Introduction

The Radio and Plasma Wave Science (RPWS) experiment is being built by an international team led by the University of Iowa for the Cassini spacecraft. This experiment will study a wide range of plasma and radio wave phenomena in the magnetosphere of Saturn and will also make scientifically important measurements during the cruise phase and at other encounters. A particular feature of the data from wave receivers is that they have a potentially vastly greater volume than the spacecraft telemetry link and onboard data handling systems are able to handle and transmit to Earth. Thus event selection, data selection and onboard data compression techniques are important for such instruments. Historically data selection has been based on hardware signal processing but recently the use of onboard software has been considered important<sup>1,2</sup>. The RPWS instrument has one processor dedicated to data compression tasks. In this paper we briefly outline the scientific data requirements for RPWS, the RPWS instrument hardware including the data compression processor (DCP) and potential DCP software structure. We then present some results of data compression tests and finally discuss the present planning for the implementation of data compression in the RPWS instrument. We note, in particular, that the complexity in the number of RPWS modes will impact on data compression yet the priority for compression will be directed at the producers of the largest data volumes.

## 2. The Scientific Requirements

Figure 1 shows the typical signal levels for the various noises which can be expected around Saturn<sup>3</sup>. A simple calculation would show that to measure all of these fully one would need - (3 antennae) x (2 x 10<sup>8</sup> samples) x 16 bps or approximately 10 Gbps for the electric components (assuming 100 MHz bandpass), and (3 antennae) x (2 x 10<sup>5</sup> samples) x 16 bps or approximately 10 Mbps for the magnetic components (assuming 100 kHz bandpass) of the wave-field. This calculation assumes that the dynamic range can be adequately quantized using a 16 bit word length. This naive and simplistic calculation does not include the telemetry allowance for a Langmuir probe or a sounder both of which are in the RPWS instrument. Typical data rates available for RPWS are a few kbps, i.e. a factor of some 10<sup>6</sup> lower than could be used. Thus the need for data compression.

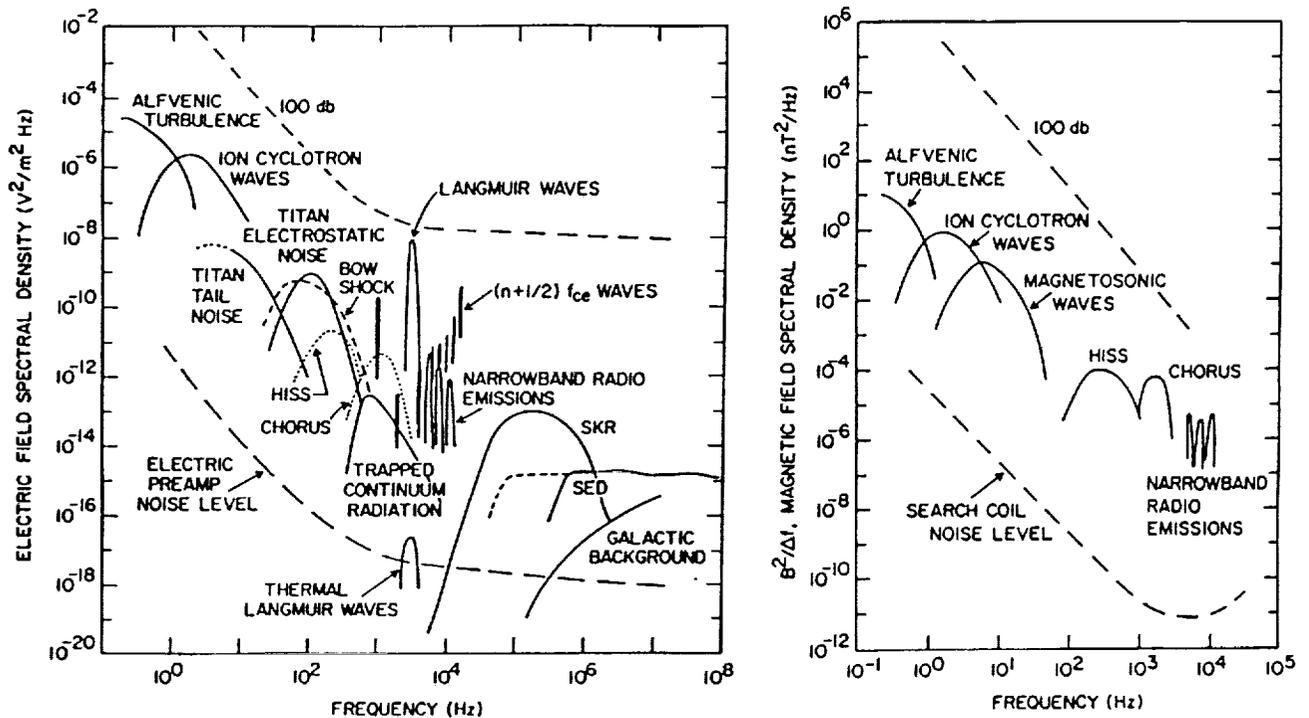


Figure 1 — The spectral characteristics of the various emissions expected to be detected during the Cassini mission to Saturn. Without any compression or selection, nearly 10<sup>8</sup> bps are required to return this information, which is a much larger value than the telemetry allows.

Clearly the RPWS team would be unwise (and unlikely to get support to try) to build sensors with such a capacity to over-produce data. In the next section we describe the RPWS instrument block diagram.

### 3. The RPWS Instrument and Data Compression Processor

Figure 2 is a simplified block diagram showing the current RPWS configuration. Three electric field antennas (configured as either a dipole plus a monopole or three monopoles), three orthogonal magnetic search coils and a Langmuir probe are the sensors. Two electric antennas can be connected to an active plasma sounder. The main signal processing blocks are the high frequency receiver (HFR), the medium frequency receiver (MFR), the five channel waveform receiver (5CWF), a digital wideband receiver (WB) and a Langmuir probe (LP). Both the HFR and MFR contain internal averaging and intensity compression circuitry, thereby reducing the onboard data handling requirement. In contrast, the 5CWF and WB receivers sample very fast, each near their respective Nyquist frequencies. Considering the WB receiver, as many as 160k samples per second are obtained in some modes. The data rates available for the RPWS investigation is around 2kbps in normal operation, but can be increased to greater than 100 kbs at predetermined times.

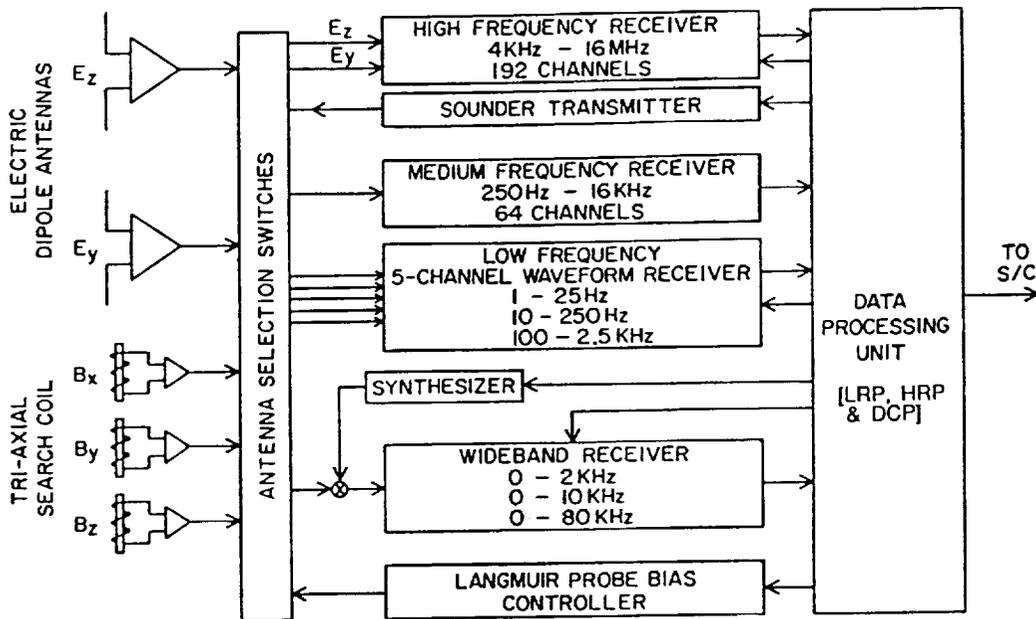


Figure 2 — A block diagram of the Cassini Radio and Plasma Wave (RPWS) experiment. Note that the Data Processing Unit (DPU) accepts signals from five different receivers. Due to the very different types of data, compression systems cannot be generic, but must be specifically tailored for each receiver.

Outputs from the signal processing blocks are taken to either the high rate data processing unit (HRP) or the low rate data processing unit (LRP). These processors are part of the instrument block labeled as the "Data Processing Unit" in Figure 2. A more detailed diagram of the DPU layout is shown in Figure 3. The HFR and MFR data will be processed by the low rate processor. In contrast, the faster-sampling 5CWF, LP, and WB data will be processed by the high rate processor. The high rate data processing unit contains a dedicated compression chip (produced by JPL using Rice's split-sample scheme) which can compress the data by a factor of approximately 2. This chip will be primarily for data from the digital wideband receiver. Other forms of data compression and data selection can be performed in the data compression processor (DCP). The DCP is connected by a single bus to both the high and low rate data processing units. This bus can only handle communications between two of the three processors at any particular time. The DCP is an 80C85 processor with 2k bytes of PROM and 64k bytes of RAM. It also includes a 16 x 16 multiply and accumulator device (Marconi MAR 7010). The DCP will perform data compression on the outputs from the various receivers, in scenarios that will be described below. Compressed data from the DCP is sent to the low rate processor where it is packetized and returned. Ultimately the spacecraft data interface is with BIU (bus interface unit) which is connected to the low rate processor.

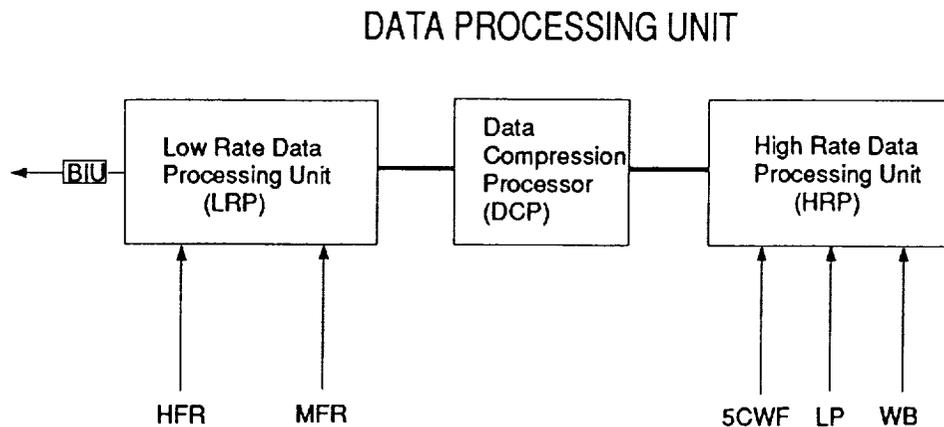


Figure 3 — A specific block diagram of the DPU. Note that internal to the HRP is a lossless compression chip. Besides this capability, there is a dedicated data compressor (DCP) that performs both selection and compression tasks.

#### 4. Algorithm Tests on Possible Compressors

In this section we present the results of some simulations and tests of data compression on various sources of data. These include tests of compression upon data obtained from a ground-based lightning sferic radio detection system, along with simulations of compression using Voyager radio wave data.

For the detection of VLF signals associated with lightning discharges, a ground-based waveform capture system was recently developed at GSFC. The system examined waveforms obtained in the frequency range between 1–30 kHz, capturing the most active interval via a “smart” selection process. Figure 4 shows a captured VLF waveform generated by a cloud-to-ground return stroke occurring on 28 August 1992. During this particular day, storms associated with the remanent of Hurricane Andrew passed near the observation site (i.e., about 50 km). The “smart” selection process identified this particular interval of time as “active” and saved the corresponding receiver output in memory. The information is further compressed using an adaptive quantization algorithm, the results of which are shown in the middle panel. In this compression scheme, the original 16-bit words are requantized to 4-bit words using 16 quantization steps equally-spaced between the minimum and maximum waveform values. This algorithm is very quick, yielding moderate, synchronous compression. Although the compression, by nature, is lossy, there is little loss of essential information concerning the lightning-generated waveform. The bottom panel of the figure displays a simple model of a typical VLF waveform from cloud-to-ground return strokes<sup>4</sup> for comparison.

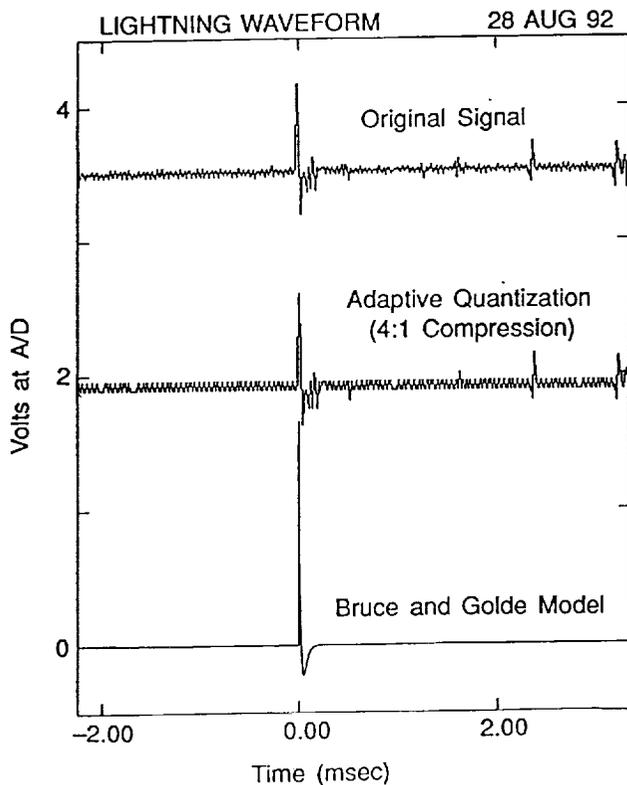


Figure 4 — A captured VLF waveform from a groundbased receiver that features a data compression system. Shown is the emission generated by a cloud-to-ground return stroke occurring on 28 August 1992, observed between 1–30 kHz. A compressed version of the waveform using adaptive quantization is also shown. Both waveforms compare well with modeled results.

There are other methods besides selection and adaptive quantization for the compression of data in the form of time series. We now present a couple examples of these methods using the data from the Voyager 2 encounter with the planet Neptune. Specifically, the data presented is from the planetary radio astronomy (PRA) experiment onboard the spacecraft. This experiment consists of a sweep frequency receiver operating between 1.2 kHz and 40 MHz<sup>5</sup>. The data from the closest approach period on 25 August 1989 is used in the study.

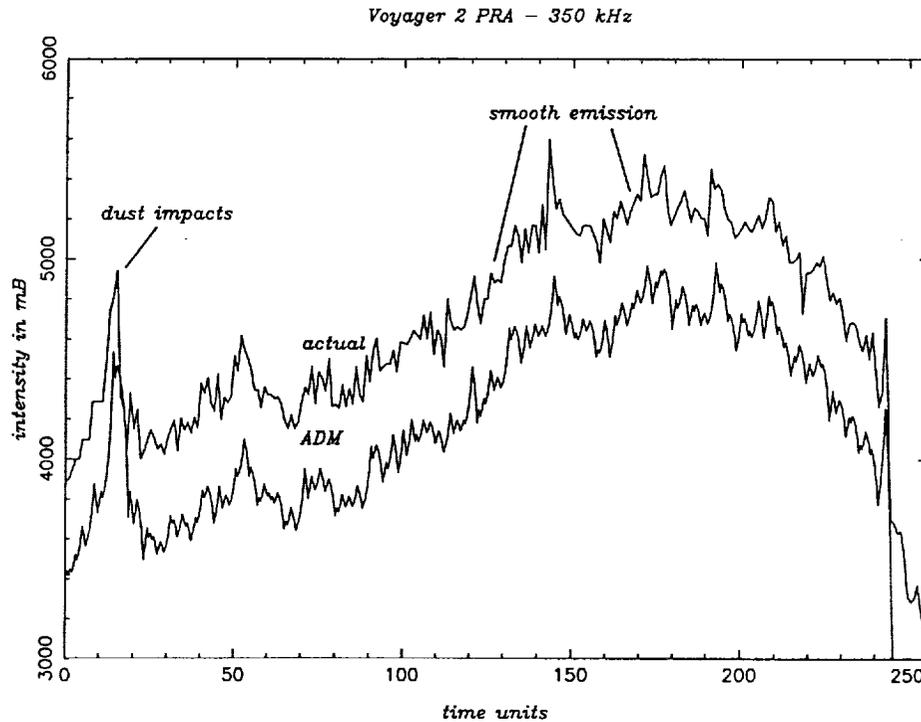


Figure 5 — A time series from the Voyager planetary radio astronomy (PRA) experiment (top curve) and its reconstruction following its compression using a 2-bit per sample adaptive delta modulation system (bottom curve).

Figure 5 shows the PRA measurements obtained during a 25 minute period just prior to closest approach to the planet. During this period, a time-averaged signal associated with dust impacts was detected by the receiver, along with a smooth emission that persisted for many hours. In the figure, the original data is presented in the top curve. The bottom curve is the same time series reconstructed following the application of an adaptive delta modulation (ADM) system commonly used in speech compression<sup>6</sup>. The ADM technique is relatively fast and yields a synchronous, 2-bit per sample output. Note that there is a reasonably good correspondence

between the actual and ADM time series. Specifically, all the signals of scientific relevance, such as the dust impacts and smooth emission, are captured by this system.

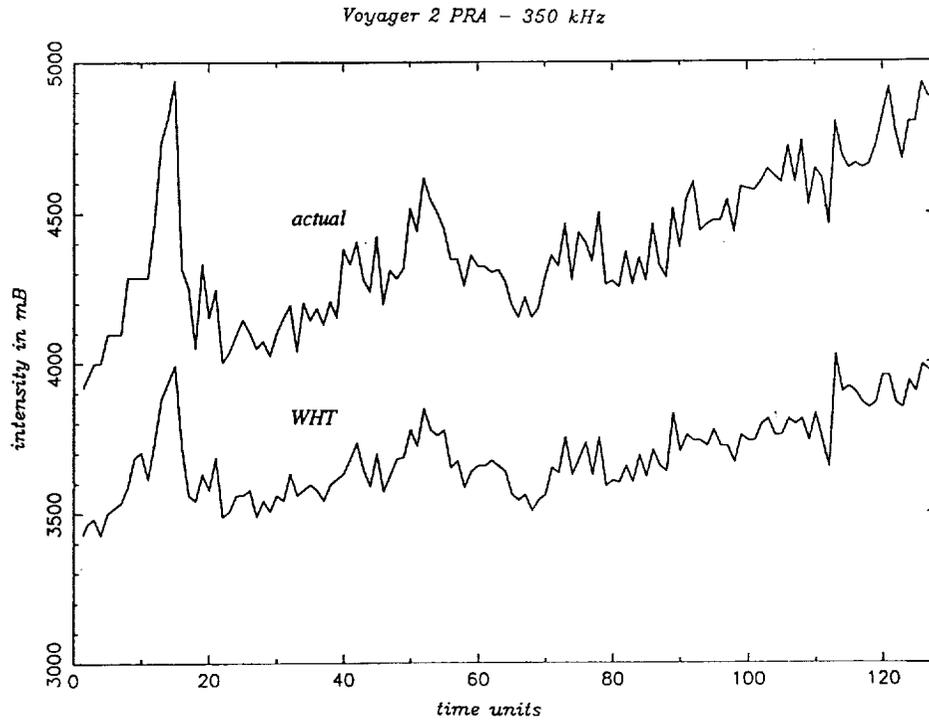


Figure 6 — A time series from the Voyager planetary radio astronomy (PRA) experiment (top curve) and its reconstruction following its compression using a Walsh-Hadamard transform/coefficient selection process (bottom curve).

In the event that there are very low data rates and available CPU time, then transform coding may be a possible means of compressing the data. Figure 6 shows a time series from the Voyager PRA experiment (top curve) along with reconstructed signal using the Walsh-Hadamard transform (WHT). The transform, itself, does not compress the data. However, once in transform space, selection of the most intense coefficients is performed. After selection, the coefficients are requantized to three bits and returned using run length encoding techniques for efficient packing in the telemetry stream. The reconstructed time series following these processes is shown in the bottom curve of Figure 6. This transform system results in compression down to about 2 bits per samples. As evident in the figure, there are some minor distortions in the WHT time series that result from the requantization process. However, the relevant scientific information is returned.

## 5. Compression Plans for RPWS

The implementation of data compression in RPWS is not yet fully decided. It is clear that the DCP will not have sufficient capacity in either processing power or memory for an extended suite of compression algorithms. Neither do we consider that it is possible to employ a single algorithm for all formats of input data, our experience suggests that algorithms need to be selected according to the nature of the data being considered. Thus, testing on the ground prior to flight will have to be performed to find the proper algorithms suitable for the different data types.

An initial approach to the Cassini RPWS compression is to compress the data from the source which produces the greatest volume of data. For example, compression and selection of events associated with the fast-sampling WB system is of primary importance due to the large data volume created by this instrument. When high rate telemetry modes ( $> 100$  kbs) are available, the data can be returned directly or undergo fast lossless compression. The dedicated compression chip produced by JPL should accomplish the latter task. However, when data rates are low ( $< 10$  kbs), data selection and lossy compression are required. One possible scenario is a selection and simple requantization process similar to that associated with Figure 4.

The DCP will also process the 5CWF and LP data. Like the WB experiment, these data sets consist of waveform measurements, but with much lower temporal resolution. Since the information rate is lower for these two receivers, the use of CPU-intensive compression systems may be possible. For example, it is desirable to transform the output from the 5CWF receiver in order to obtain spectral information. Compression in transform space may then be possible.

The HFR and MFR return spectral information averaged over a predefined time interval (usually averaged values in 10's of milliseconds). The data points are correlated in both frequency and time, thus a number of compression systems are possible. However, the compression applied is dependent upon the data sets usage. If radio spectrograms of limited resolution are created as the final output, lossy compression with large compression ratios may be used. Adaptive delta modulation is one possible system. Since the operation times of these receivers are relatively slow, transform coding and coefficient selection may also be possible within the constraints of the processor. However, if measurements are going to be crosscorrelated, exact values are needed, and these can only be realized using lossless compression. As described below, the final usage of the data product on the ground is another driver in the selection of the compression system.

## 6. Discussion

The discussion of the use of data compression for scientific data raises the question of what, if

any, degradation in the quality of the data is acceptable. For wave data where the data in question are a time series of field amplitude measurements it may be appropriate to guarantee that the additional uncertainty, or noise, in the data is below the noise level of the receiver, typically by around 3dB or so lower. For frequency spectra the scientifically important information are usually plotted on a grey-scale or colour scaled plot with a restricted quantization of the number of levels. For these data, then, a data compression to a word length somewhat longer than that needed to code the number of colour levels will usually be entirely adequate. When, for example, polarization and directions of the incoming wave are being measured then it is important not to have any data degradation. In general, the data compression tests by simulations should show that no scientific result is changed or be an artefact of the compression strategy which was selected.

It is evident that the choice of data compression system depends upon the receiver and its associated data product, DCP processing time, and the final data usage on the ground. Choosing the correct compression systems for each receiver will be difficult, but might be performed with the use of logic diagrams, like that shown in Figure 7. Illustrated is the possible compression scenario for the MFR data product. Note that very different compression systems could be applied, depending upon the data's final usage. In situation where there are multiple paths to the same final visualization product, one path will be selected based upon the compressibility and speed of the system. For example, to create frequency versus time spectrograms, adaptive delta modulation, transform coding, or adaptive quantizing are just a few of the possible scenarios, given an infinite amount of processing strength. In reality, the RPWS system uses an 8085 processor, thus only the most simple compression is possible. Thus, adaptive quantizing or delta modulation might may be selected in lieu of a transform coding system.

### COMPRESSION SCENARIO: MFR

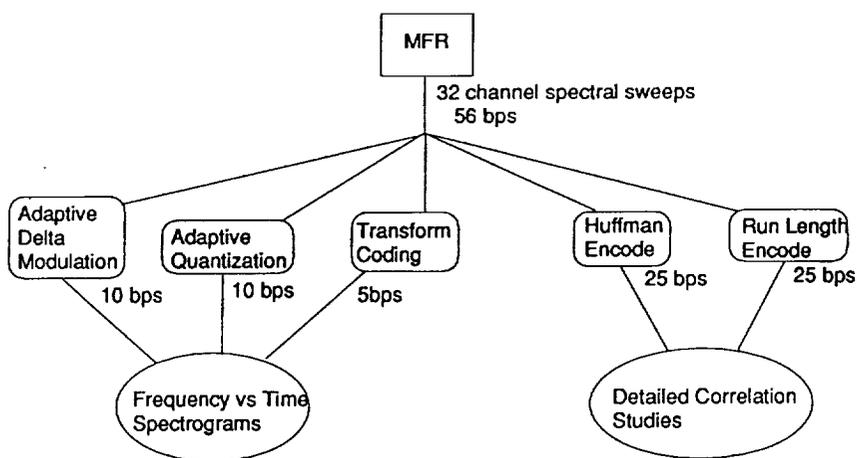


Figure 7 — An illustration of the logic that will be implemented to select the data compression systems for the DCP.

## Conclusions

This brief report gives an overview of the Cassini RPWS data processing system. At this time, the hardware portion is well defined and is selected based upon weight, power, and development costs considerations. The corresponding software portion is more complicated, with many different scenarios possible depending upon the application of the data on the ground. Future work includes algorithm selection, testing, and development for the DCP.

## References

- [1] Gough M.P. and L.J.C. Woolliscroft, Microprocessors in Space Instrumentation, *Space Technology*, 9, 305 - 313, 1989.
- [2] Woolliscroft, L.J.C., et al., Cassini radio and plasma wave investigation: Data compression and scientific applications, submitted, *J. Brit. Interplanetary Soc.*, 1992.
- [3] Gurnett, D. A., *A plasma and radio wave science investigation for the Cassini orbiter*, Proposal to NASA, 1990.
- [4] Bruce, C. E. R., and R. H. Golde, The lightning discharge, *J. Inst. Elec. Eng.*, 88, 487-505, 1941.
- [5] Warwick, J. W. et al., Planetary radio astronomy experiment for Voyager missions, *Space Sci. Rev.*, 21, 309-327, 1977.
- [6] Song, C.-L. et al., A variable-step robust delta modulator, *IEEE Trans. Comm. Tech.*, COM-19, 1033-1044, 1971.



REPORT DOCUMENTATION PAGE			Form Approved OMB No. 0704-0188	
Public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188), Washington, DC 20503.				
1. AGENCY USE ONLY (Leave blank)	2. REPORT DATE April 1993	3. REPORT TYPE AND DATES COVERED Conference Publication		
4. TITLE AND SUBTITLE  1993 Space and Earth Science Data Compression Workshop			5. FUNDING NUMBERS  936	
6. AUTHOR(S)  James C. Tilton, Editor				
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)  Goddard Space Flight Center Greenbelt, Maryland 20771			8. PERFORMING ORGANIZATION  93B00034	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)  National Aeronautics and Space Administration Washington, D.C. 20546-0001			10. SPONSORING/MONITORING AGENCY REPORT NUMBER  NASA CP-3191	
11. SUPPLEMENTARY NOTES This workshop was organized and sponsored by NASA, and cosponsored by the IEEE Computer Society Technical Committee on Computer Communications (TCCC). It was held in conjunction with the 1993 Data Compression Conference (DCC '93), which was cosponsored by the IEEE TCCC in cooperation with NASA and the Center of Excellence in Space and Information Sciences, Goddard Space Flight Center.				
12a. DISTRIBUTION/AVAILABILITY STATEMENT  Unclassified-Unlimited Subject Category 59			12b. DISTRIBUTION CODE	
13. ABSTRACT (Maximum 200 words)  This document is the proceedings from a "Space and Earth Science Data Compression Workshop," which was held on April 2, 1993, at the Snowbird Conference Center in Snowbird, Utah. This workshop was held in conjunction with the 1993 Data Compression Conference (DCC '93), which was held at the same location, March 30 to April 1, 1993. The workshop explored opportunities for data compression to enhance the collection and analysis of space and Earth science data. The workshop consisted of eleven papers presented in four sessions. These papers described research that is integrated into, or has the potential of being integrated into, a particular space and/or Earth science data information system. Presenters were encouraged to take into account the scientist's data requirements, and the constraints imposed by the data collection, transmission, distribution, and archival system. The workshop was organized by James C. Tilton of the NASA Goddard Space Flight Center, Sam Dolinar of the Jet Propulsion Laboratory, Sherry Chuang of the NASA Ames Research Center, and Dan Glover of the NASA Lewis Research Center.				
14. SUBJECT TERMS  Data Compression, Image Compression, Signal Processing, Image Processing, Space Science, Earth Science			15. NUMBER OF PAGES 128	
			16. PRICE CODE A07	
17. SECURITY CLASSIFICATION OF REPORT Unclassified	18. SECURITY CLASSIFICATION OF THIS PAGE Unclassified	19. SECURITY CLASSIFICATION OF ABSTRACT Unclassified	20. LIMITATION OF ABSTRACT Unlimited	