

A Reliable Multicast for XTP

Bert J. Dempsey
Alfred C. Weaver

JOHNSON
GRANT
IN-62-CR

Digital Technology

August, 1990

115099
P.10

N92-34108

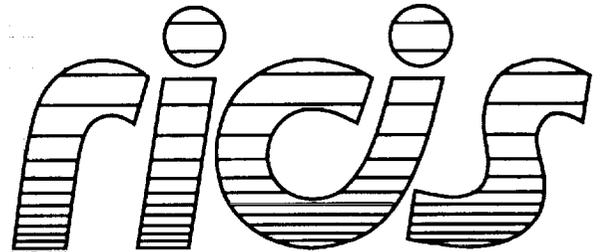
Unclas

63/62 0115099

(NASA-CR-190632) A RELIABLE
MULTICAST FOR XTP Interim Report
(Research Inst. for Computing and
Information Systems) 10 p

Cooperative Agreement NCC 9-16
Research Activity No. SE.31

NASA Johnson Space Center
Engineering Directorate
Flight Data Systems Division



Research Institute for Computing and Information Systems
University of Houston-Clear Lake

INTERIM REPORT

The RICIS Concept

The University of Houston-Clear Lake established the Research Institute for Computing and Information Systems (RICIS) in 1986 to encourage the NASA Johnson Space Center (JSC) and local industry to actively support research in the computing and information sciences. As part of this endeavor, UHCL proposed a partnership with JSC to jointly define and manage an integrated program of research in advanced data processing technology needed for JSC's main missions, including administrative, engineering and science responsibilities. JSC agreed and entered into a continuing cooperative agreement with UHCL beginning in May 1986, to jointly plan and execute such research through RICIS. Additionally, under Cooperative Agreement NCC 9-16, computing and educational facilities are shared by the two institutions to conduct the research.

The UHCL/RICIS mission is to conduct, coordinate, and disseminate research and professional level education in computing and information systems to serve the needs of the government, industry, community and academia. RICIS combines resources of UHCL and its gateway affiliates to research and develop materials, prototypes and publications on topics of mutual interest to its sponsors and researchers. Within UHCL, the mission is being implemented through interdisciplinary involvement of faculty and students from each of the four schools: Business and Public Administration, Education, Human Sciences and Humanities, and Natural and Applied Sciences. RICIS also collaborates with industry in a companion program. This program is focused on serving the research and advanced development needs of industry.

Moreover, UHCL established relationships with other universities and research organizations, having common research interests, to provide additional sources of expertise to conduct needed research. For example, UHCL has entered into a special partnership with Texas A&M University to help oversee RICIS research and education programs, while other research organizations are involved via the "gateway" concept.

A major role of RICIS then is to find the best match of sponsors, researchers and research objectives to advance knowledge in the computing and information sciences. RICIS, working jointly with its sponsors, advises on research needs, recommends principals for conducting the research, provides technical and administrative support to coordinate the research and integrates technical results into the goals of UHCL, NASA/JSC and industry.

RICIS Preface

This research was conducted under auspices of the Research Institute for Computing and Information Systems by Bert J. Dempsey and Alfred C. Weaver of Digital Technology. Dr. George Collins, Associate Professor of Computer Systems Design, served as RICIS research coordinator.

Funding was provided by the Engineering Directorate, NASA/JSC through Cooperative Agreement NCC 9-16 between the NASA Johnson Space Center and the University of Houston-Clear Lake. The NASA research coordinator for this activity was Frank W. Miller of the Systems Development Branch, Flight Data Systems Division, Engineering Directorate, NASA/JSC.

The views and conclusions contained in this report are those of the authors and should not be interpreted as representative of the official policies, either express or implied, of UHCL, RICIS, NASA or the United States Government.



A Reliable Multicast for XTP

Bert J. Dempsey and Alfred C. Weaver

Department of Computer Science
Thornton Hall
University of Virginia
Charlottesville, Virginia 22903
(804) 924-7605
bjd7p@virginia.edu, weaver@virginia.edu

1. Types of Multicast Service

Multicast services needed for current distributed applications on LANs fall generally into one of three categories: *datagram*, *semi-reliable*, and *reliable* (Figure 1). Transport layer multicast datagrams represent unreliable service in which the transmitting context 'fires and forgets'. XTP executes these semantics when the MULTI and NOERR mode bits are both set. Distributing sensor data and other applications in which application-level error recovery strategies are appropriate benefit from the efficiency in multideestination delivery offered by *datagram* service. *Semi-reliable* service refers to multicasting in which the control algorithms of the transport layer — error, flow, and rate control — are used in transferring the multicast distribution to the set of receiving contexts, the *multicast group*. The multicast defined in XTP provides *semi-reliable* service. Since, under a *semi-reliable* service, joining a multicast group means listening on the group address and entails no coordination with other members, a *semi-reliable* facility can be used for communication between a client and a server group as well as true peer-to-peer group communication. Resource location in a LAN is an important application domain. The term 'semi-reliable' refers to the fact that group membership changes go undetected. No attempt is made to assess the current membership of the group at any time — before, during, or after — the data transfer.

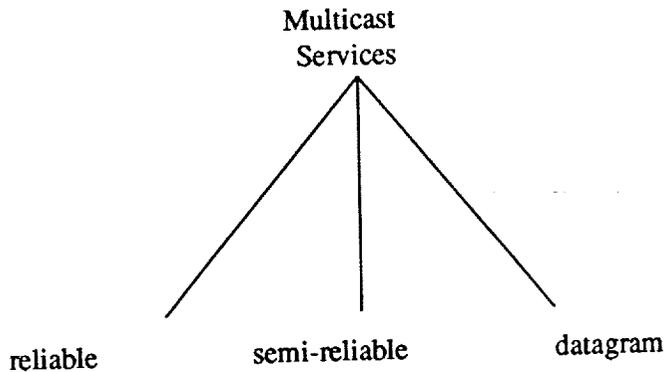


Figure 1 — Multicast Services

Applications that require a strong degree of data and behavior consistency between the process group members such as distributed databases, replicated program execution modules, and automated control programs need *reliable* multicast service. These applications use peer-to-peer communication to update global group state or data copies. The desired multicast semantics include (1) delivery of each message to either all or none of the operational processes, (2) receipt of messages in the same order by all group members, and (3) detection of group membership changes and reporting of these changes as they occurred relative to the message delivery stream.

Various research efforts address the need for *reliable* multicast. Message ordering and reliability can be achieved by having a single group member ([2]) or special system node ([4]) that assigns a sequence number to all messages to the group. The sequencer takes on the responsibility of holding past messages in a *history buffer* and retransmitting to group members that detect lost messages through discontinuities in the sequence of message numbers. Tseung ([6]) uses a systems hardware solution by adding special network nodes that perform the sequencing and logging tasks, thus addressing connectivity errors as well. In [5] a token circulates between group members. The token-based scheme permits all of the requirements of

reliable multicast service to be met as well as preserving atomic delivery semantics during network partitions and consistent message sequencing when process groups overlap. Finally, Birman and Joseph [1] [3] explore a range of ordering properties for distributed applications and propose a two-phase commit protocol for the FIFO ordering needed by *reliable* multicast.

2. Name Space Management

All group communication schemes need a Directory Services entity to manage the name space for multicast groups. Providing this support in a distributed fashion can be accomplished in the following fashion. The local name manager entity receives a group creation call from the user and hashes into the set of available addresses for an address to associate with the new group. The manager then broadcasts a request to the well-known address of all name managers and waits for any reply indicating a collision. This broadcast is performed some number of times for reliability before the local manager decides that the address can be safely assigned to the new group. The group creator becomes the first member of the new group, and the new group address is broadcast to the set of local managers. For real-time systems where the long latency of group creation may be unacceptable, a (human) user-managed table of 'permanent groups' should be available to shortcircuit run-time verification of address uniqueness.

3. A Reliable Multicast Service Layered onto XTP

XTP can provide datagram and semi-reliable multicast services. The full functionality for a reliable multicast service, however, clearly lies outside the realm of a transport layer protocol. Below is an outline of how a reliable multicast service can be layered onto XTP. Using the extra channel available through tagged data (i.e. data carried in BTAGs and ETAGs) a higher layer *reliable multicast protocol* (RMP) can be constructed. The hook into XTP that would enable this protocol is a control bit associated with tagged data fields indicating that these fields are carrying RMP information.

Figure 2 shows one of the four members of a multicast group transmitting to the group. The RMP entity at the transmitter manages reliability for the one-to-many (forward direction) data flow via control information solicited from and sent by remote RMP entities. If the number of remote listeners is known at the transmitter (which RMP guarantees), then the local instance of RMP sets up and manages that number of unicast connections for reverse direction data flow. Each listener establishes a reverse direction connection, which is an ordinary XTP unicast connection. In both the forward and reverse directions, RMP entities multiplex user data and RMP control information (in the tagged data channel) over a single connection.

Message serialization under RMP is trivial since only one member can be transmitting messages at any time. Since each transmission requires the participation of all group members, the n-way connection set-up described here serves as an exclusive lock on the group. The state of message delivery to the group is available in the progress of the current transmitter's one-to-many connection. This eliminates the need for an explicit global ordering of messages to the group and simplifies the task of identifying messages during failure recovery.

RMP can achieve atomic delivery of messages by performing a two-phase commit. The RMP transmitter notifies RMP group members to withhold a message from the user until all users send confirmation, through tagged data in the reverse connections, that they have the message. Then the transmitter issues a 'commit' directive. Failure recovery techniques must ensure that if any group member receives the 'commit' and delivers the message that all operational members will deliver the message as well.

Once a group is created, a new member joins (or leaves) the group by contacting the current transmitter in the group. The transmitter uses a two-phase commit mechanism to atomically change the ordered list of group members kept by each member. Notification of the

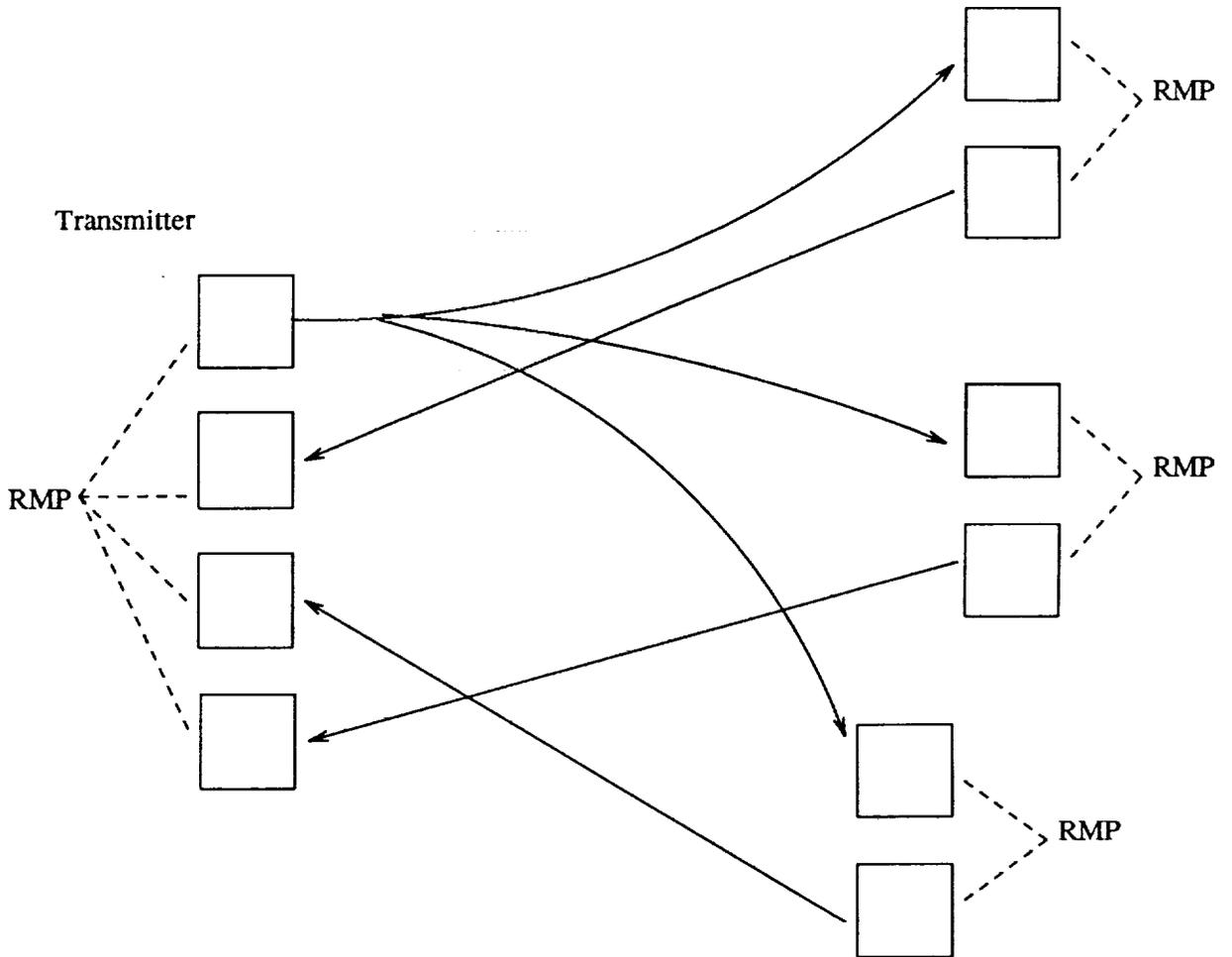


Figure 2 — Reliable Multicast with Reverse Connections

change will be passed up to the application, which will act on the information in an application-dependent fashion. By coordinating membership changes on a message boundary, the changes can be sequenced in the message stream being delivered to each local user (Figure 3). The transmitter brings the joining member into (or drops the leaving member out of) the conversation at the correct message boundary. Like user-generated messages, membership change notifications will be delivered to either none or all operational members in the same

order. If no member of the group is currently transmitting, the new (or leaving) member can itself contact the group and coordinate its joining (or leaving).

Members leaving the group due to process or host failures will be detected by the current (or next) transmitter to the group when the transfer fails due to the reception of less than the proper number of control messages. This action triggers failure recovery mechanisms to rebuild the ordered list of members, clean up any state inconsistencies (e.g. messages delivered to only a subset of the operational process group members), and recalculate delivery parameters affected by the group membership change (e.g., rate control and timer defaults). These recovery actions can be handled by a token-passing mechanism, similar to that in [5], using reliable unicasts. If no secondary failures occur, the token will need to circulate twice around the virtual

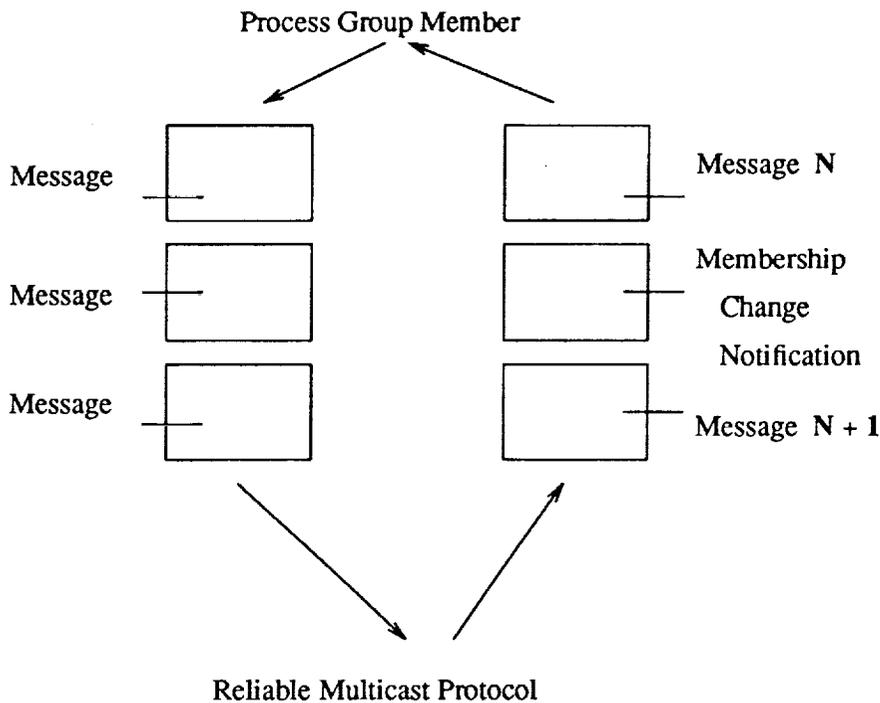


Figure 3 — User Notification on Group Membership Change

ring of group members in order to create and commit the new list at each member site.

References

1. K. Birman and T. Joseph, Reliable Communication in the Presence of Failures, *ACM Transactions on Computer Systems* 5,1 (February 1987), 47-76.
2. J. Chang and N. F. Maxemchuk, Reliable Broadcast Protocols, *ACM Transactions on Computer Science* 2,3 (Aug. 1984), 251-273.
3. T. Joseph and K. Birman, Reliable Broadcast Protocols, in *Distributed Systems*, S. Mullender (editor), ACM Press, 1989, 293-319.
4. M. F. Kaashoek, A. S. Tanenbaum, S. F. Hummel and H. E. Bal, An Efficient Reliable Broadcast Protocol, *Operating Systems Review* 23,4 (October 1989).
5. B. Rajagopalan and P. McKinley, A Token-Based Protocol for Reliable, Ordered Multicast Communication, *Proceedings of Eighth Symposium on Reliable Distributed Systems*, Seattle, Washington, October 1989.
6. L. C. N. Tseung, Guaranteed, Reliable, Secure Broadcast Networks, *IEEE Network*, November 1987, 33-37.