

# A Centralized, Tree-Based Approach to Network Repair Service for Multicast Streaming Media

Dan Rubenstein  
Computer Science Department  
UMass Amherst  
drubenst@cs.umass.edu

Nicholas F. Maxemchuk  
AT&T Labs - Research  
Florham Park, N.J.  
nfm@research.att.com

David Shur  
AT&T Labs - Research  
Middletown, N.J.  
shur@research.att.com

## Abstract

*IP multicast provides best-effort delivery. Packets encounter variable delays and may be lost because of transmission errors and buffer overflows. Real-time multimedia streaming services require that most packets arrive at the receivers prior to an application deadline. Multicast quality on the current Internet is often inadequate for these applications. We have solved this problem by placing repair servers inside the network. The repair servers recover missing packets by communicating with each other, then re-multicast the repaired stream to nearby receivers on a new address. Multicast reception in the constrained area is typically much better than in the wide area Internet.*

*In this paper we address the problem of constructing a repair graph. The repair graph shows which repair servers each repair server communicates with to recover missing messages. Our objectives when constructing this graph conflict with each other. We want high reliability: every repair server to recover as many missing messages as possible, as quickly as possible. But we also want low cost: this recovery should use as little of the network bandwidth as possible. We present a centralized algorithm to generate repair graphs. We demonstrate through simulation that these graphs achieve a level of reliability that is almost as high as that achieved by repair graphs specifically designed for high reliability. At the same time, our graphs maintain a cost that is almost as low as the cost in repair graphs designed for low cost.*

## 1 Introduction

The Internet's flexible design enables it to transmit a variety of types of traffic. Traffic can be as simple as a point-to-point data communication, or as complex as a large-scale, multimedia collaboration that simultaneously transmits information between numerous hosts. The deployment of IP Multicast [1] significantly reduces the load imposed on the Internet by large-scale, multi-host communications. However, IP Multicast is best-effort: a packet transmitted via IP Multicast need not reach all intended receivers. Hence, on

its own, IP Multicast is unsuitable for sessions that require low packet loss rates.

Multimedia broadcast sessions, such as live or pre-recorded broadcast radio / TV, are applications that, via IP Multicast, can be made widely available to consumers at a low cost. However, peoples' interest in tuning in to such sessions decreases with the quality of the reception. Since packet loss degrades reception quality, it is necessary to add some additional mechanism on top of the best-effort IP Multicast service to keep loss rates at low levels.

Here, we present a system that improves reception quality by recovering from packet losses in streams transmitted via IP Multicast. The system utilizes a few dozen servers inside the network as *repair servers*. Each server caches packets from the original transmission, and can forward (via unicast) packets to other repair servers that did not receive the packets via the original transmission. After a fixed delay on the order of several seconds, a repair server then multicasts the packets in scoped sub-regions of the network, providing receivers within the scoped region with a transmission that is delayed by a few seconds relative to the original multicast session. However, this delayed transmission contains fewer losses, and thus presents a higher quality copy of the original transmission than what the receiver would have obtained had it joined the original transmission group.

We focus on producing a recovery system that can easily be deployed by a network service provider or by a set of cooperating network service providers. The system has several unique features. First, it is built to support existing and future applications that are designed to interact with the current best-effort IP multicast model. Our system does not require that these applications possess any additional functionality to handle late packet arrivals or to request repairs. Second, the system is fault tolerant: it uses a simple algorithm to route repair transmissions around any repair server that ceases to operate, and reconfigures itself to avoid subsequent failures. Last, the system's configuration information is maintained at a central point of control, simplifying system operation, monitoring, and on-line debugging.

There has been a significant amount of prior work that uses repair servers in a similar manner [2, 3, 4, 5, 6]. What

separates our work from prior art is our focus on keeping the architecture simple where it doesn't need to be complicated. This simplicity of the architecture is captured in terms of three basic design principles:

1. Decisions are centralized, as long as the centralization does not compromise the scalability of the system, and as long as the system can operate correctly and remain productive for an extended time if separated from the centralized decision point. In particular, we use a centralized algorithm to determine from which repair servers a given repair server should request repairs. Other works use distributed algorithms that use various probing techniques to identify appropriate repair servers [2, 3, 5, 6, 7], which complicates their deployment.
2. The complexity of the protocol at the repair servers is kept simple, and the amount of state that a repair server must maintain is kept low.
3. The system utilizes a simple model of the underlying network. A system designed to utilize a complex network model tends to be quite complex itself, and makes it difficult to assess its performance, or debug problems that may arise.

This simplification increases the likelihood of a successful deployment, and also simplifies the task of monitoring and debugging.

This paper specifically focuses on the protocol that is utilized by the repair system to repair packet losses at the repair servers, allowing them to deliver high quality multimedia multicast sessions to receivers. Our protocol provides high *reliability*: a high likelihood that a repair server can deliver a packet to a receiver within the receiver application's deadline. This is accomplished using a small amount of additional bandwidth on network links, keeping the operating cost low (in terms of link usage). In the protocol, repair servers obtain missing packets from other repair servers. What affects the reliability and cost of the protocol is the choice of *repair graph*: a graph whose nodes are the repair servers, and whose edges indicate the directions in which repairs flow between repair servers. We demonstrate how a variant on the minimum spanning tree algorithm generates repair graphs that can achieve high reliability at a low cost. The only knowledge of the underlying network needed by these algorithms is the Euclidean distance between all pairs of repair servers, and between each repair server and the originating multicast source. We demonstrate through simulation that, given the limited knowledge we have about the underlying network, the repair graph generated by the algorithm achieves a level of reliability that is almost as high as that achieved by repair graphs specifically designed for high reliability (repair servers request repairs from servers that lie near the transmitting source). At the same time, it maintains

a cost that is almost as low as the cost maintained by the repair graphs designed to maintain a low cost (repair servers request repairs from nearby servers).

The paper proceeds as follows. In Section 2, we introduce the network model: our abstract view of the network and of the capabilities of the repair servers for which we wish to provide repair service. Section 3 presents the algorithms we use to generate the repair tree, and motivates why we believe these algorithms produce the kinds of repair graphs that would be most effective at providing high reliability at a low cost. Section 4 presents performance results of simulations of our repair graph algorithm on the network model. We discuss related work in Section 5, and last we discuss future directions and conclude in Section 6.

## 2 Network Model

In this section, we present our abstract model of the network on top of which we will be building a repair graph. We begin by describing the application model that we support, and then describe our abstraction of the Internet topology that we believe is a reasonable basis on which we should design our repair graph algorithm.

### 2.1 Application Model

Our repair graph algorithm is designed to connect repair servers that support multimedia broadcast applications, where the applications are designed to run on top of best-effort IP multicast. As a result, the repair server cannot rely on the application to actively participate in the repair service protocol (e.g., repair requests, playout buffering). Our model of the interface between repair server and application is based on the model presented in [8], which we now summarize.

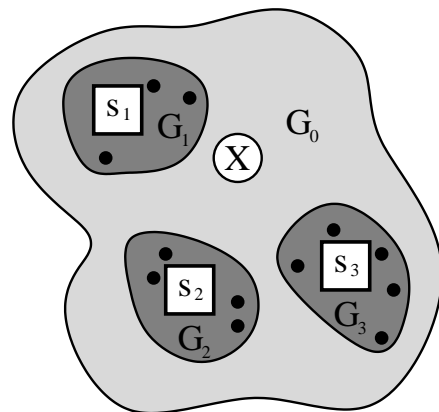


Figure 1: The Non-modified Application Repair Server Architecture

Figure 1 depicts how multicast groups are used by the repair system to improve reception quality at receivers. The

transmission source, represented by an 'X', transmits its signal on a multicast group,  $G_0$ . Transmissions from the source onto  $G_0$  have a large scope (e.g., a ttl of 127), and can be received in the lightly-shaded region. Some of the data generated by the source is lost within the network as it propagates to receivers (black dots) or repair servers (square white boxes) joined to multicast group  $G_0$ . Hence, receivers that join to  $G_0$  will often receive a poor quality signal due to substantial packet losses.

Each repair server,  $s_i$ , joins to group  $G_0$ , and delays its transmission of the received data by some fixed amount of time,  $t_i$ , beyond its scheduled playback time. During this period of time, it detects packet losses in its received stream on group  $G_0$ , and attempts to recover these lost packets from neighboring repair servers via unicast requests. The recovery attempt consists of a series of one or more unicast requests to nearby repair servers, which unicast back the requested packet if it is available. When a time of  $t_i$  has elapsed beyond the playback time of a received or recovered packet, the packet is transmitted by the repair server,  $s_i$ , on a separate multicast group,  $G_i$ . Transmissions to  $G_i$  by  $s_i$  are scoped to within a local region (e.g., ttl of 0,1 or slightly larger). In Figure 1, we indicate the smaller scoping of the group  $G_i$  by a darker shaded region which surrounds the repair server.

Rather than join group  $G_0$ , a receiver that lies within the transmission range of some repair server,  $s_i$ , can join to group  $G_i$  and receive a delayed, but much higher quality transmission. This is because a receiver near a repair server  $s_i$  that is joined to group  $G_i$  receives almost all the packets that were received or recovered by  $s_i$ . If the repair graph is configured in a reasonable manner, then it is highly likely that a repair server  $s_i$  will be able to recover a packet as long as some repair server upstream on the graph receives the packet in the original transmission.<sup>1</sup> The repair will propagate downstream due to the chain of unicast requests and retransmissions of the lost packet along each hop of the repair graph. Thus, the likelihood is high that a receiver will receive a packet on  $G_i$  that it would have lost on  $G_0$ .

Note that because the source and receiver do not participate in the recovery process, applications that have been and are being developed to use the standard IP multicast UDP-like interface can utilize the protocol by simply joining a repair multicast group instead of joining the original multicast group. Further details can be found in [8].

Repair servers do not maintain state about who has requested a repair. If a repair server,  $s_i$ , receives a request for a packet from a repair server,  $s_j$ , and it cannot provide the packet,  $s_i$  makes no response, and keeps no information regarding the request. Hence, it is  $s_j$ 's responsibility to query  $s_i$  at some future point in time if it still desires to receive the packet from  $s_i$ . If each repair server were to maintain state about queries, it would be possible to propagate a repair through a chain of servers that lost the initial transmission in much less time. However, maintaining such state

<sup>1</sup>A node  $A$  is *upstream* from another node  $B$  if there is a directed path from  $A$  to  $B$ .

increases the complexity of the repair server code (which requires a re-request mechanism in any event because requests themselves might be lost), and is not necessary for the range of applications that we wish to support with our protocol.

In our system, the organization of the repair graph is formulated from a centralized process, which then contacts each repair server to inform it of its connection information. This centralized process assigns a set of parents to each repair server from which the server may request packet repairs. One can construct the repair graph from this information by drawing edges from each repair server to its set of parents. The centralized process can contact servers on occasions when the repair graph needs to be updated (e.g., a repair server in the system fails). However, we expect such updates to be infrequent. We have already built a repair system that operates in this manner that locates and routes around faults in the network. Furthermore, repair servers whose services are not needed (no receivers are joined to the repair session, and no repair servers are requesting repairs) can enter a sleeping mode that conserves resources at the host running the repair server process, as well as in the neighboring network region. The description of these attributes of the system is discussed in [9].

## 2.2 Abstraction of the Underlying Multicast Network

We evaluate our system through simulation over an abstract model of the multicast network. Our abstraction of the underlying multicast network attempts to capture several realistic features of Internet multicast routing, and at the same time, avoid unnecessary details that we expect will not significantly affect protocol performance. Since we compare the performance of repair graphs whose nodes consist of the transmission source and repair servers, our underlying network graph in the model only contains nodes that can potentially serve as repair servers. We do this rather than introduce additional complexity in the model by including intermediate nodes that connect repair servers. We expect that removing these additional nodes does not significantly alter our results. Our model also does not include the last hop to the receivers that connect to the repair servers. Because one of our requirements is that application-level code cannot be modified to request or detect repair packets, there is little that can be done to improve the communication performance between a receiver and its "nearest repair server". We leave it to the receiving application (or user of the application) to its "nearest repair server", i.e., to choose a group on which it obtains the best service it can, given the limitations on the number of repair servers that can be deployed.

### 2.2.1 Underlying topology

We now discuss our construction of the underlying network topology. Network nodes, each of which represents a potential source or repair server, are placed on a two-dimensional

grid in which only a fixed subset of grid squares can contain nodes. The grid squares which contain nodes represent populated areas, such as cities, or, if one prefers a larger scale, continents. The grid squares that do not contain nodes represent unpopulated regions or oceans. We use a scaled down version of what is used in [10] to generate our graphs. In [10], up to 10,000 receivers are used per sample. Here, we are only interested in the placement of repair servers, and expect on the order of 50 nodes to be sufficient for providing a global repair service: an intelligent placement of these servers throughout the network will improve the transmission quality for a large majority of receivers within the network. We randomly select 5 grid squares from a 5x5 square grid to be populated regions. The remaining squares are oceans (devoid of nodes). We refer to each square that contains nodes as a *continent*.

After placing the nodes within the network on the continents, we randomly assign bi-directional links to connect the nodes, where the probability of constructing a given link is a function that decreases with the distance between the nodes it is to connect. Thus, nodes that are close together are more likely to be directly connected (it follows that the density of connectivity is higher within a continent than across continents). The algorithm to generate links does not terminate until a path exists between all pairs of nodes.

### 2.2.2 Building the multicast tree

Once we have generated the underlying network topology, we choose one node to be the source of the multicast group, and 20 of the remaining 49 unused nodes to be repair servers that will participate in the session. The multicast tree is the shortest-path tree from the source to each of the repair servers, where a path can proceed through any node within the underlying graph (whether or not the node has been selected as a repair server). We could perhaps build trees that are more realistic by increasing the aggregate number of nodes in the graph beyond 50. Doing so would likely increase the expected hop-count between repair servers. However, we expect that ISPs will choose repair servers in a strategic manner, and the shortest path between two nearby repair servers is likely to closely approximate the Euclidean distance between them.

### 2.2.3 Loss and recovery

We represent packet loss on the multicast tree as a Bernoulli process on each link of the tree. A repair server fails to receive a packet whenever any link upstream from it (toward the source) drops a packet. This process captures spatial, but not temporal, loss correlation observed for multicast transmission. In this paper, we will examine the system’s ability to recover a single packet that is lost in parts of the network. For this reason, the fact that we do not capture the temporal loss correlations between successive transmissions is irrelevant.

For the purposes of building the repair graph, we permit each repair server to connect directly to any other repair server (i.e., the graph in which the 20 repair servers are the nodes and the edges represent paths for direct communication is a 20-clique). Again we make the assumption that the distance along the communication path between two repair servers closely approximates the Euclidean distance. Constructing the graph as such also implies that it is quite possible for two servers to have a more direct mode for communication than the path that connects them within the multicast tree.

We are interested in testing the performance of our system under two types of losses on top of the repair graph. First, we assume that the combined request and repair transmissions between two repair servers failing to deliver the repair from the requestee to the requester is a Bernoulli loss process. Second, we assume that it is possible that a small number of repair servers (0, 1, or 2) may fail during the session. These servers do not respond at all to requests for repairs.

### 2.2.4 Repair server algorithm

A repair server that detects a missing packet requests a repair immediately, and continues to do so periodically until the packet is received, or the packet’s deadline for playout on the locally scoped multicast group expires. In the current implementation, the period of the request is a second. Since all repair servers are likely to detect a packet loss at approximately the same time (relative to a second), the  $i$ th request from all servers that have not received the packet occur at approximately the same time, and hence it makes sense to model the request process as a series of *rounds*. In each round, a repair server that has not yet recovered the packet makes a request to a parent for the packet. Because repair servers do not maintain state for packet requests, a repair propagates at most one hop on the repair graph per round toward a repair server that needs the packet.

## 3 A Repair-Graph Building Algorithm

We now describe the algorithm we use to generate repair graphs. Let us quickly review the traits of a “good” repair graph. We define the *reliability* of a repair graph more formally as the expected number of packets, averaged over all repair servers, that can be recovered before the deadline. We want a repair graph that exhibits a high reliability, close to 1. At the same time, we want to limit our usage of network resources. Here we assume that the cost of a transmission between repair servers equals the Euclidean distance between those servers. This is not an unreasonable assumption. At present, leased-line providers charge by the mile.

We also impose several additional requirements based on observations of trials of our prototype implementation [9]. The algorithm should be robust in environments where there

are losses of repair packet transmissions, as well as over single node failures. Robustness over multiple node failures is of course preferred. However, it is unlikely that two nodes' times of failure will overlap.<sup>2</sup> Last, we assume that the only information that is available to the algorithm is the geographical location of the source and of the repair servers in the network: this information allows us to compute the Euclidean distance between pairs of repair servers and between the source and any repair server. Thus, we will not need to obtain accurate routing information, nor do we require initial estimates of the loss rates between various pairs of repair servers.

### 3.1 Toward a Robust Repair Graph

Before embarking on a description of our repair-graph building algorithm, we first describe at a higher level what motivated us to construct the algorithm the way we did.

#### 3.1.1 A ring at the top

If the data source for a session participates in the repair protocol, then ensuring that a repair server eventually receives a copy of a lost packet (barring a node failure, a flushing of the source's cache, or a deadline expiring) is simply a matter of rooting the repair graph at the source (making the source upstream from all repair servers). However, our assumption that the application code cannot be modified does not allow us to include the data source in the repair graph. This means that it is conceivable that in some instances, the original transmission of a packet is lost by all repair servers, and is therefore unrecoverable. While we cannot hope to recover from such occurrences, we can try to construct a repair graph that gives a high likelihood of recovery in a small number of rounds in most loss scenarios by ensuring that a few hops upstream, there is a repair server that has a high likelihood of receiving the packet.

One method for locating such a repair server would be to observe the traffic for some period of time, and then perform a correlation study for each repair server to determine which server would best suit its needs for repair. However, such a process complicates the protocol, and with traffic patterns changing over time in the network, it is not clear that such a solution would be effective. Instead, we make the following observation: nodes with a smaller Euclidean distance to the source are more likely to lie upstream (nearer to the source) on the underlying multicast tree, and it follows that in most cases, it is more likely that these servers will receive the original data transmission, and subsequently be able to provide repairs.

Our solution is to identify a small set of repair servers that are near the source, such that it is most likely that if any

---

<sup>2</sup>The central point of control maintains periodic contact with each server at a rate such that the system is capable of detecting node failures in under a minute [9], and in this same period of time it should be possible to restructure the repair graph to avoid the failed node.

repair server receives a packet, then the packet was received by some repair server in this small set. We then connect the servers in the small set in the form of a ring so that if any of them receive the packet, then all of them should receive the packet as the packet propagates around the ring. The remaining repair servers are connected to the ring from downstream, so that any packet that is received by some repair server within the ring can subsequently be obtained by any repair servers downstream.

We note that if it is possible to co-locate a repair server at the location of the data source, then the process of locating a ring of repair servers that surround the source can be omitted. For the remainder of the paper, we assume that the source is located at a point that does not permit the co-location of a repair server, necessitating the construction of the ring.

#### 3.1.2 Two Trees

A "good" repair graph efficiently distributes repairs from servers that are more likely to receive data in the initial transmission to those servers that are less likely to receive the data. Most often, a tree is constructed for this purpose [3, 4, 5, 7], since it is the most efficient means of producing a fully connected graph (a path exists from the root of the tree to each node). However, a node failure, unless lying at a leaf of the tree, would partition the tree and prevent propagation of repairs from nodes on one side of the partition to the other. To overcome this drawback, we need to construct a repair graph that has the efficiency close to that of a tree, but also has the ability to route around such node failures.

### 3.2 Our solution

At a high level, our solution is the following:

- Select the 3 nodes that are likely to surround the source and construct a ring from these nodes with edges that is bi-directional (directed edges are chosen going both in the clockwise and counter-clockwise direction).
- Treating the three nodes in the ring as a single node, construct a primary tree with high reliability and low cost from the remaining nodes, rooted at the node that represents the ring.
- Remove the (directed) edges used in the primary tree, and generate a secondary tree rooted at the ring on the remaining edges. We would like to construct this secondary tree such that the graph built by concatenating the ring and two trees together results in a graph that can route around any single node failure.
- During each round, a repair server that has not yet received a copy of the packet chooses an edge on the graph generated by the concatenation of the ring and the two trees, and requests the repair from the node at the other end of the edge.

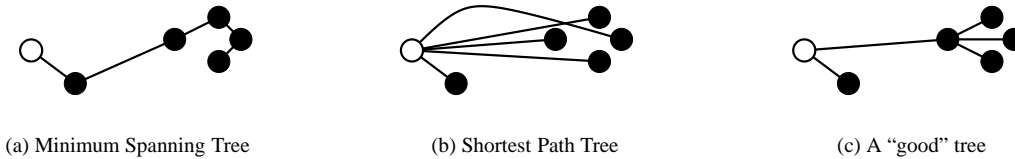


Figure 2: Various Trees built on top of a 6-clique, where the white node is the root.

We now describe the algorithm to construct the trees. We postpone our discussion of how we choose the three nodes from which we build the ring. For the time-being, the reader should assume that the three node ring has already been formed. The resulting tree that is constructed by the following algorithm extends from the ring, i.e., the ring can be thought of as a super-node that roots the tree.

The algorithm that constructs the primary tree is a modification of an algorithm that generates a minimum-spanning tree (MST). An MST is the tree with the lowest link cost. However, there are some nodes in an MST whose path from the root contain many hops. In our model, since packets are dropped on each hop with equal probability, and since packet repairs traverse a single hop in each round for a fixed number of rounds, the level of reliability decreases at a repair server as the number of hops to repair servers increases. To increase reliability, the number of hops from the ring to repair servers must be decreased. The tree with the fewest number of hops to each repair server is a shortest path tree rooted at the ring (where a path's length is the number of hops contained within the path). Since edges are permitted between any pair of repair servers, the shortest path tree connects each repair server directly to a repair server that lies within the ring. This introduces several long, high-cost edges into the repair graph. Ideally, we want to build a tree where long links in a repair graph are shared by several repair servers that are near to one another. For instance, looking at this in terms of a network spread out over several continents, it makes sense to form trees where only a single link connects several repair servers across continents, but the path from a repair server to this cross-continental link has a small number of hops.

Figure 2 shows the difference between a minimum spanning tree, a shortest path tree, and a tree we consider ideal. In all three graphs, the node locations are identical, only the edges used to connect the nodes change. In the minimum spanning tree (Figure 2(a)), there is a node that is 5 hops from the source (white node). In the shortest path tree (Figure 2(b)), all nodes are one hop from the source, but most are connected using long-length edges. In the "good" tree (Figure 2(c)), most edges are of short length, and nodes are at most two hops from the source.

### 3.2.1 Construction of the primary tree

The algorithm to construct the primary tree adds one node at a time, and relies on two measures of distance within the tree and network. We write  $c(n, m)$  for the Euclidean distance from a node  $n$  to a node  $m$ , and measure the network cost for sending a repair over link  $(n, m)$  as being proportional to  $c(n, m)$ . We define  $h_1(n)$  to be the minimum number of hops it takes along in the primary tree to travel from a node in the ring at the root of the tree to node  $n$ . Note that  $h_1(n)$  is only defined for a node  $n$  after the node has been added to the tree.

1. Let  $N$  be the set of nodes without a path on the tree to the ring.
2. Let  $M$  be the set of nodes with a path on the tree to the ring.
3. While  $N$  is non-empty
4. Select  $n \in N$  and  $m \in M$  such that  $c(n, m) + \alpha h_1(m) \leq c(n', m') + \alpha h_1(m')$  for all other  $n' \in N, m' \in M$ .
5. Add directed edge  $(n, m)$  to the graph.
6. end while

Figure 3: An iteration of the algorithm to build Tree 1

Figure 3 presents the main iteration within the body of the algorithm used to generate the primary tree. The three nodes that form the ring near the source are connected and placed in the set,  $M$ , of nodes which have a path on the tree to a node within the ring.  $N$  is initially the set of the remaining nodes. Each time the iteration of the body is performed, a node is moved from  $M$  to  $N$ . When  $M$  is empty, all nodes are attached to the tree and the algorithm is complete. The value chosen for the parameter,  $\alpha$ , affects the shape of the primary tree. If  $\alpha$  is zero, then the algorithm generates an MST: the link with lowest cost that connects to the tree and does not form any cycles is always added on a given iteration. As  $\alpha$  is increased, more emphasis is placed on minimizing the number of hops from the ring to a repair server. As  $\alpha$  tends toward infinity, the resulting tree converges toward the shortest path tree (where path length is the number of hops in the path).

We now describe the construction of the ring that forms

the root of the repair tree. The objective is to try and “surround” the paths on the multicast tree from the source without underlying knowledge of the multicast tree. We want to keep small the number of nodes that we place within the ring to minimize the maximum hop-count to propagate a repair to all repair servers in the ring. To accomplish these tasks, we run the algorithm to build the primary tree, rooted at the source, and then perform a breadth-first search on the resulting tree. When two nodes have the same path length from the source, we give precedence to the node with the smaller aggregate edge cost. The first three nodes we come across are the nodes that form the ring. In the ring that roots the primary tree, directed edges are chosen such that the third node obtained in the breadth first search recovers from the second node, which recovers from the first, which recovers from the third. After this ring is formed, the algorithm to generate the primary tree is rerun, with  $M$  initially set to the three nodes that form the ring. This result is a set of trees, which, when joined to the set of edges in the ring in one direction, is a single tree plus one additional edge.

### 3.2.2 Construction of the secondary tree

For the secondary tree, the ring at the root of the tree contains the same set of nodes as the ring that roots the primary tree, but the edge directions are reversed (i.e., the first node added to the ring recovers from the second, which recovers from the third, which recovers from the first). All directed edges used within the primary tree are removed from the underlying graph so that the secondary tree and primary tree share no common, directed edges. Note, however, that if directed edge  $(m, n)$  is in the primary tree, we permit the edge  $(n, m)$  in the secondary tree.

The main purpose of the secondary tree is to ensure that the graph built from the concatenation of the two trees can route around any node failure. We accomplish this by restricting the set of edges that can be used in the secondary tree to those that satisfy certain depth requirements. We define  $h_2(n)$  in a manner similar to the definition of  $h_1(n)$ : a value is assigned to  $h_2(n)$  once node  $n$  has been added to the secondary tree, where  $h_2(n)$  equals the minimum number of hops it takes along in the secondary tree to travel from a node in the ring at the root of the tree to node  $n$ .

Figure 4 presents the main iteration within the body of the algorithm used to generate the secondary tree. Note that this algorithm is similar to the algorithm used to build the primary tree: a node is added to the tree when, compared to the other nodes that can be added, it minimizes a combination of link cost and hop count. As in the algorithm that builds the primary tree, the relative weights of the link cost and hop count depend on the choice of  $\alpha$ . This algorithm imposes an additional requirement on the node being added: a node,  $n$ , can attach downstream from a node  $m$  only if  $h_1(n) \geq h_1(m)$  within the primary tree. This additional property guarantees that the graph formed from the edges from both trees can route around any single node fail-

1. Let  $N$  be the set of nodes without a path on the tree to the ring.
2. Let  $M$  be the set of nodes with a path on the tree to the ring.
3. While  $N$  is non-empty
4.   Select  $n \in N$  and  $m \in M$  where  $h_1(n) \geq h_1(m)$ , such that  $c(n, m) + \alpha h_2(m) \leq c(n', m') + \alpha h_2(m')$  for all other  $n' \in N, m' \in M$  where  $h_1(n') \geq h_1(m')$ .
5.   Add directed edge  $(n, m)$  to the graph.
6. end while

Figure 4: An iteration of the algorithm to build Tree 2

ure. This is proven formally in [11]. The proofs are omitted here due to space restrictions.

We conclude this section by noting each node in the ring has two parents (the other two nodes in the ring). Each node that is not in the ring also has two parents: a node in each tree, or, for each time its parent is the ring (represented as a super-node), it has a distinct node in the ring as a parent. It follows that if each repair server requests repairs from both of its parents, any repair that was received by a node within the ring will eventually reach all other nodes.

## 4 Evaluation

In this section, we evaluate the performance of our algorithm via simulation. Our evaluation proceeds in three steps. First, we determine an appropriate value for the tunable parameter,  $\alpha$ , that achieves a low cost while maintaining a high likelihood of packet recovery (i.e., reliability) for all receivers. We will see that using a value of  $\alpha = 1$  within the algorithm results in trees that use long, higher cost links infrequently (keeping cost low), and have a low hop-count for most repair servers to recover the packet (keeping reliability high). Next, we consider the impact of building the repair graph as the concatenation of two trees rather than building it as just a single tree. We find that utilizing the links on the secondary tree causes a small increase in cost and a small reduction in reliability in the case where there are no node failures, but that when node failures occur, requesting repairs from the secondary tree can increase reliability significantly. Last, we examine the impact on reliability and cost as we vary the loss rate on both the multicast tree and on the recovery tree. We find that the loss rate on the multicast tree reduces reliability of the repair graph, but this reduction applies to any generated repair graph. Hence, the repair graph generated with  $\alpha = 1$  in which both repair trees are utilized in the repair process remains the preferred option. We find that loss on the repair graph has much less of an impact on both cost and reliability than having this loss on the original multicast

distribution tree.

We run simulations as follows: A simulation configuration consists of a choice for  $\alpha$ , a fixed loss rate on links of the multicast tree, a fixed loss rate on links of the repair tree, a fixed number (0, 1 or 2) of nodes that fail (are unable to forward repairs), and a round schedule that indicates per round from which tree (primary or secondary) a repair server selects its parent to issue its repair request. For each configuration, we perform 500 runs. Each run consists of the following steps: First, we randomly generate a network, and randomly choose 21 out of the 50 nodes to be possible repair servers and multicast source. We then choose 10 different nodes at random to be the multicast source. For each choice of multicast source, we generate a multicast tree and repair graph (where the remaining 20 nodes chosen are repair servers). In the case where we have either one or two node failures, we iterate over all possible combinations of nodes failing, considering each combination once. We then transmit 10 packets, record how many repair servers were able to recover the packet by the end of each round, and what the cumulative cost of transmissions was at the end of each round. We average over all results, giving equal weight to each packet transmission in each run, to obtain the reliability and average cost.

Figure 5 demonstrates how varying  $\alpha$  impacts the reliability and cost of a repair session. Here, there is no loss of transmissions between repair servers, and there are no node failures. The loss rate on each link during the original multicast transmission of the packet is .03. Each repair server issues a request along the primary tree during the first four rounds, and switches to the secondary tree during the last four rounds, and gives up trying to obtain a packet after eight rounds. On the  $x$ -axis of both Figures 5(a) and 5(b), we vary the number of rounds. The  $y$ -axis of Figure 5(a) indicates the average level of reliability, and the  $y$ -axis of Figure 5(b) indicates the average cost. Points plotted at  $x = 0$  in Figure 5(a) indicate levels of reliability of the original multicast session (without any repairs). The various curves plot reliability and cost as a function of round for different values of  $\alpha$ .

We see that by increasing  $\alpha$ , we can increase the reliability of the repair graph. This is because a high value of  $\alpha$  increases the weight of the hop-count component within the metric that the tree generating algorithm tries to minimize. As a result, the length of the path from a node in the ring at the root of the graph to the repair server is decreased. This in turn decreases both spatial correlation of loss and the maximum number of hops that a repair needs to travel in worst case scenarios. However, cost increases, since repair servers are less likely to choose parents that are nearby neighbors, and more likely to choose parents that are near to the source. A value of  $\alpha = 10$  is most often a star, with each repair server connecting directly to some node in the ring at the root. A value of  $\alpha = 0$  produces a minimum spanning tree (if we consider the ring at the root as a single node before applying the MST algorithm).

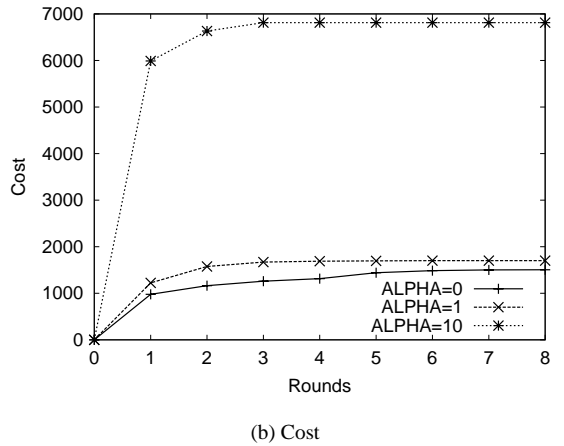
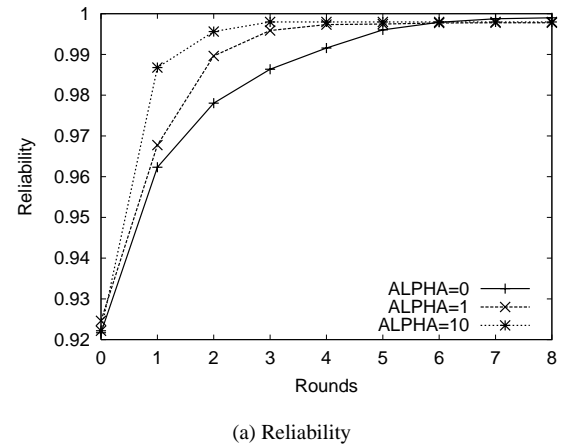
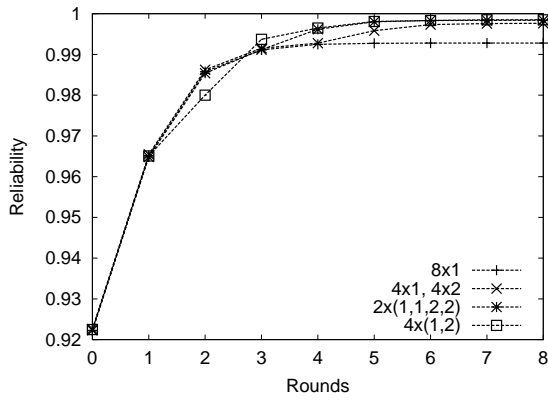


Figure 5: How changing  $\alpha$  changes reliability and cost

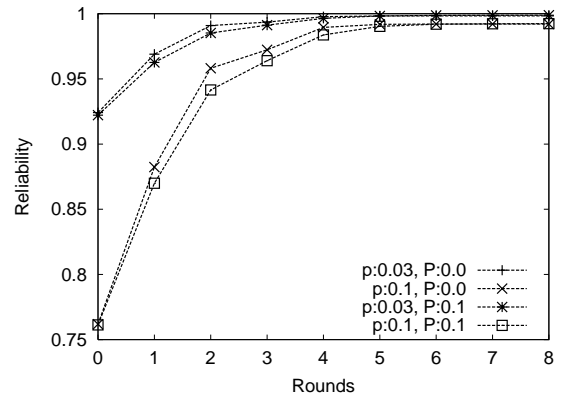
We see a significant increase in reliability as we vary  $\alpha$  from 0 to 1, and a significant increase as well when it is varied from 1 to 10. However, the increase in cost as we vary  $\alpha$  from 0 to 1 is not nearly as significant as it is when it is varied from 1 to 10. We conclude that it makes the most sense to choose a value of  $\alpha$  close to 1, since this gives a significantly higher level of reliability than that for  $\alpha = 0$  without a significant increase in cost.

Figure 6 demonstrates how utilizing a second tree can increase reliability in the event of node failures. Because our simulations use identical loss rates on all links, one does not achieve any significant gain in reliability by utilizing a second graph when there are no node failures. However, using a single repair tree, a node failure might partition a set of repair servers that lost a packet from the set of nodes that might be able to repair the loss. Figure 6(a) demonstrates the impact that a single node failure has (recall we average over all possible configurations of the location of the node failure), and Figure 6(b) demonstrates the impact of two node

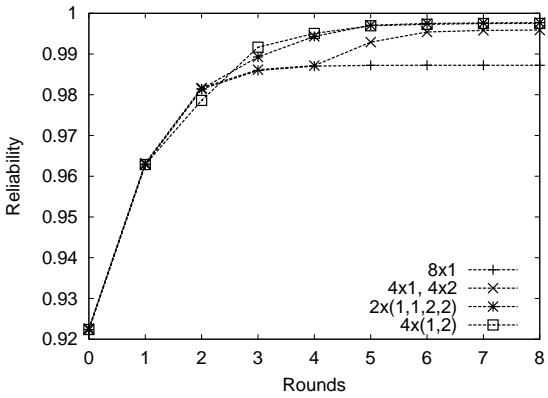




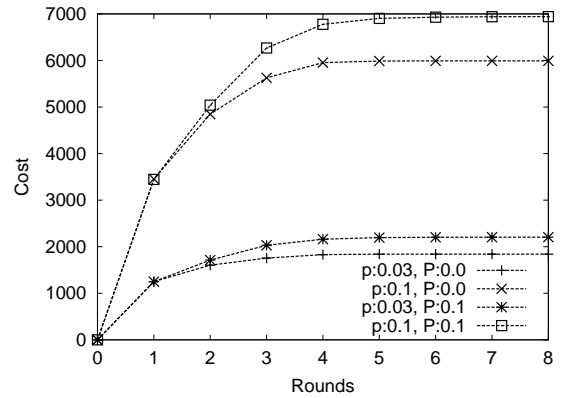
(a) One node failure



(a) Reliability



(b) Two node failures



(b) Cost

Figure 6: The second tree’s impact on node failures

Figure 7: The impact of loss on reliability and cost

failures at once. Again we vary the number of rounds on the  $x$ -axis and the reliability on the  $y$ -axis. The different curves represent different orderings of repair graphs used per round. For instance, the curve labeled “4x(1,2)” means that on odd rounds, the parent is chosen from the primary tree, and from the secondary tree on even rounds. We see that switching to the secondary tree allows a higher level of reliability, and this is exemplified in the unlikely event that there is more than a single node failure at any given point in time.

Figure 7 demonstrates the effect of varying the loss rate on the original multicast tree (labeled  $p$ ) and on the repair tree (labeled  $P$ ). We see that varying the loss rate on the original multicast tree has a larger impact on both the reliability and cost of the repair server system than similar variations in loss rate on the repair transmissions. This is not surprising, since increasing the loss rate on the original multicast tree increases the expected number of receivers that need additional repairs. Increasing the loss rate on the repair tree only affects the reliability of the subset of receivers that

did not receive the initial transmission, and therefore this has a smaller impact on overall reliability.

Last, we find that the conclusions drawn by examination of Figures 5 and 6 remain the same for the loss rates depicted in Figure 7. In other words, we find that for reasonable loss rates, a value of  $\alpha = 1$  and alternating parents between the two repair trees over the various rounds provides a repair system that exhibits the most desirable tradeoff between low cost and high reliability.

## 5 Related Work

IP multicast [1] is designed as a best-effort service, shifting the responsibility of providing some sort of reliable delivery to the protocol layers that lie above the network layer. Several early works examine reliable multicast in networks where no explicit knowledge of participant location was necessary [12, 13]. However, the current preferred means to

achieve a reliable multicast that scales to a large number of receivers is to form some sort of hierarchical structure in the network. The idea was first developed into a protocol by Paul et al [3] by building the hierarchy from participating receivers.

Providing better-than-best-effort service for deadline-driven traffic, such as audio and video applications, was examined independently by Maxemchuk et al in [8], Xu et al in [2], and Lucas et al in [5]. In these works, receivers request repairs from other receivers. The goal is to repair as many lost packets before their deadlines as possible, rather than providing full reliability, as would be necessary for a data-transfer. Xu's work was extended in [6] to demonstrate the performance benefits of including *waypoints*: servers inside the network dedicated to improving protocol performance. In the application presented in that paper, waypoints act as repair servers, and reduce the repairing load on receivers in the session.

A variety of works [2, 6, 7, 14, 15] present different algorithms for generating the repair hierarchy that can be used to provide reliable multicast. All of these approaches build their repair graphs dynamically using distributed algorithms. A dynamic algorithm to build the repair graph can customize the graph precisely to the current networking conditions. The distributed nature allows the algorithm to scale to a large number of tree participants. However, our results indicate that a simple, static, centralized algorithm builds a repair graph that is sufficient for our needs, both in terms of reliability and scalability.

## 6 Future Directions and Conclusion

We have proposed and examined a simple, static, centralized algorithm to build a repair graph to provide resilient multicast support to real-time multicast sessions that can tolerate small playback delays. The simplicity of the algorithm is due to the high level of abstraction of our underlying network model, and because decisions that are made infrequently are centralized. We start by limiting the set of features in the network that we model to only those that we believe are important, and design an algorithm that provides a high level of reliability while maintaining a low link utilization cost for this simple model. The simplicity of the algorithm makes the system easy to deploy, and its static, centralized nature facilitates understanding and/or debugging of the protocol in a real Internet environment.

We are currently incorporating our algorithm into our existing repair service system and will next examine its effectiveness at providing repairs within the actual Internet. It would also be of interest to compare the level of reliability and link utilization cost of our approach with the protocols that build their repair graphs in a distributed, dynamic fashion to see whether the performance gains due to dynamic adaptation warrant the resulting additional complications in deployment.

## References

- [1] S. Deering and D. Cheriton. Multicasting routing in datagram internetworks and extended LANs. *ACM Trans. on Computer Systems*, 8(2):85–110, May 1990.
- [2] R.X. Xu, A.C. Meyers, and H. Zhang. Resilient Multicast Support for Continuous Media Applications. In *Proceedings of NOSSDAV'97*, St. Louis, MO, July 1997.
- [3] S. Paul, K.K. Sabnani, J.C. Lin, and S. Bhattacharyya. Reliable Multicast Transport Protocol (RMTP). *IEEE JSAC*, 15(3):407–421, April 1997.
- [4] S. Kasera, J. Kurose, and D. Towsley. A Comparison of Server-Based and Receiver-Based Local Recovery Approaches for Scalable Reliable Multicast. In *Proceedings of INFOCOM'98*, San Francisco, CA, March 1998.
- [5] M.T. Lucas, B.J. Dempsey, and A.C. Weaver. MESH: Distributed Error Recovery for Multimedia Streams in Wide-Area Multicast. In *Proceedings of IC3N'97*, 1997.
- [6] K. Sripanidkulchai, A.C. Meyers, and H. Zhang. A Third-Party Value-Added Network Service Approach to Reliable Multicast. In *Proceedings of ACM SIGMETRICS'99*, Atlanta, GA, May 1999.
- [7] B. Levine, D. Lavo, and J.J. Garcia-Luna-Aceves. The Case for Concurrent Reliable Multicasting Using Shared Ack Trees. In *Proceedings of ACM Multimedia'96*, Boston, MA, November 1996.
- [8] N.F. Maxemchuk, K. Padmanabhan, and S. Lo. A Cooperative Packet Recovery Protocol for Multicast Video. In *Proceedings of ICNP'97*, Atlanta, GA, October 1997.
- [9] D. Rubenstein, D. Shur, and N.F. Maxemchuk. Repair Server Activation Implementation. *Technical Report in preparation*, December 1999.
- [10] N.F. Maxemchuk. Video Distribution on Multicast Networks. *IEEE Journal on Selected Areas in Communications*, 15(3):357–372, April 1997.
- [11] D. Rubenstein, N.F. Maxemchuk, and D. Shur. A Centralized Approach to Network Repair Service for Multicast Streaming Media. *AT&T Technical Memorandum TM HAI7200000-991129-03*, September 1999.
- [12] D. Towsley, J. Kurose, and S. Pingali. A Comparison of Sender-Initiated and Receiver-Initiated Reliable Multicast Protocols. *IEEE JSAC*, 15(3):398–406, April 1997.
- [13] S. Floyd, V. Jacobson, C. Liu, S. McCanne, and L. Zhang. A Reliable Multicast Framework for Light-Weight Sessions and Application Level Framing. *IEEE/ACM Transactions on Networking*, 5(6):784–803, December 1997.
- [14] M. Hofmann and M. Rohrmuller. Impact of Virtual Group Structure on Multicast Performance. In *Proceedings of 4th COST237 Workshop*, Barcelona, Spain, December 1997.
- [15] B. N. Levine, S. Paul, and J.J. Garcia-Luna-Aceves. Organizing multicast receivers deterministically according to packet-loss correlation. In *Proceedings of ACM Multimedia'98*, Bristol, U.K., September 1998.