

The Impact of Multicast Layering on Network Fairness*

Dan Rubenstein, Jim Kurose, and Don Towsley

<http://www-net.cs.umass.edu/{~drubenst, ~kurose, ~towsley}>

Technical Report 99-08
Department of Computer Science
February, 1999

A shorter version of this report will appear in ACM SIGCOMM '99.

Abstract

Many definitions of fairness for multicast networks assume that sessions are single-rate, requiring that each multicast session transmits data to all of its receivers at the same rate. These definitions do not account for multi-rate approaches, such as layering, that permit receiving rates within a session to be chosen independently. We identify four desirable fairness properties for multicast networks, derived from properties that hold within the max-min fair allocations of unicast networks. We extend the definition of multicast max-min fairness to networks that contain multi-rate sessions, and show that all four fairness properties hold in a multi-rate max-min fair allocation, but need not hold in a single-rate max-min fair allocation. We then show that multi-rate max-min fair rate allocations can be achieved via intra-session coordinated joins and leaves of multicast groups. However, in the absence of coordination, the resulting max-min fair rate allocation uses link bandwidth inefficiently, and does not exhibit some of the desirable fairness properties. We evaluate this inefficiency for several layered multi-rate congestion control schemes, and find that, in a protocol where the sender coordinates joins, this inefficiency has minimal impact on desirable fairness properties. Our results indicate that sender-coordinated layered protocols show promise for achieving desirable fairness properties for allocations in large-scale multicast networks.

1 Introduction

The current Internet has few internal mechanisms to regulate the rates at which sessions should transmit data. How to achieve *fairness* within such a network, in effect allowing sessions to share bandwidth in a manner that satisfies some set of network utilization criteria, remains a challenging research problem. The problem is further complicated in networks that support both unicast and multicast delivery services. Current definitions of multicast fairness [17, 13, 3, 19, 6] typically assume that sessions are *single-rate*, requiring all receivers within a multicast session to receive data at a uniform rate. However, *layered multicast* permits *multi-rate* transmission: different receivers within a session can receive data at different rates. This is accomplished by layering data among several multicast groups and allowing each receiver to independently determine the subset of layers (i.e., multicast groups) it joins. Protocols have used a layered approach to support multicast applications ranging from live multimedia [11, 8, 9, 1] to reliable data transfer [14, 18, 4]. These protocols have the appealing property that the transmission rate to each receiver is constrained only by the bandwidth availability on the receiver's own data-path from the data source, and is not limited by other receivers' rate limitations in the same session. The fairness literature does suggest intuitions about how layering might increase the set of desirable fairness properties that hold for a particular fair allocation of receiver rates. What is lacking, however, is a formal study that examines the impact that layering has on fair allocations within a large-scale multicast network.

The goal of this paper is to contribute to the formal understanding of how layering impacts fairness in multicast networks. In particular, we focus on how layering affects properties of multicast max-min fairness in an environment

*This material was supported in part by the National Science Foundation under Grant No. NCR-9508274, NCR-9527163, CDA-9502639, and by DARPA under Grant No. N66001-97-C-8513. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

in which each session has a single sender. We have chosen to use max-min fairness as our fairness measure since its formal definition is a well-accepted criterion for fairness, allowing us to proceed directly to an examination of the properties of a fair allocation. We believe that with other definitions of fairness, layered approaches will yield similar fairness advantages, and expect this work to stimulate interest in examining the impact of layering for these other definitions.

We show that allowing multicast sessions to be multi-rate instead of single-rate “improves” max-min fairness within a network. To do this, we identify four desirable fairness properties of a max-min fair allocation of unicast sessions. One simple example of a property is that receiver rates should be equal for two receivers whose data transmission paths from their respective senders traverse an identical set of links. We examine multicast max-min fair allocations under the definition given by Tzeng and Siu [17], that requires that all sessions are single-rate, and find that several of these fairness properties do not necessarily hold within the max-min fair allocation (the two receiver rate example presented above is one such property). We extend the multicast max-min fair definition to permit multi-rate sessions, and formally prove that, when all sessions in a network are multi-rate, all of our identified fairness properties hold for the max-min fair allocation. We also consider networks in which not all sessions are multi-rate (e.g., a session may have an application-specific requirement that requires it to be single-rate), and examine the effect on fairness properties of the max-min fair allocation as single-rate sessions are “replaced” by identical multi-rate sessions (i.e., same session members, same topology). We demonstrate, using our identified set of fairness properties and a mathematical ordering relation of allocations that indicates an allocation’s “level” of max-min fairness, that increasing the set of “replaced” sessions results in an increase in the “level” of max-min fairness and that more fairness properties hold for max-min fair allocations.

Next, we consider several practical limitations encountered by sessions that use layering to achieve max-min fair rates. We show that if each receiver’s fair rate is restricted to what can be obtained by joining some fixed set of layers, a max-min fair allocation need not even exist. However, we do demonstrate that receivers can achieve an *average* rate that matches their fair rate by using precisely timed joins and leaves. These joins and leaves must be tightly coordinated among receivers in the same session (i.e., correlating their sets of received packets) in order to prevent excess bandwidth utilization on a shared link. We introduce the notion of *redundancy*, a ratio of bandwidth used in practice by a session on a shared link to the theoretical lower bound needed on that link to deliver fair rates to downstream receivers, in order to quantify bandwidth usage. While several works have identified the negative implications of redundancy, there has not yet been an analytical exploration of its magnitude, or of its effect on fair allocations within a network. We show that increased redundancy leads to a decrease in the “level” of max-min fairness, to a decrease in the number of fairness properties that hold for the max-min fair allocation, and, usually, to a decrease in receivers’ fair rates. We examine how the ideas in [11, 8, 18], that coordinate joins of receivers within a session, significantly reduce the negative effects of redundancy. The examination is performed via analytical modeling and simulation of max-min fair congestion control protocols in which receivers join and leave layers based on congestion observations. Within the model, we present three protocols that differ in the degree to which the layer joins are coordinated among session receivers. We find that although redundancy is still not optimal, coordinated joins reduce redundancy most significantly when the correlation in loss among receivers is high, and that a protocol with sender coordination keeps redundancy at low enough levels to allow layered multicast to achieve non-bandwidth-wasteful fairness within a multi-rate multicast network.

The paper proceeds as follows. Section 2 presents theoretical results for multicast max-min fairness with multi-rate sessions. Section 3 introduces the notion of redundancy, and Section 4 examines the effects of join coordination in several simple congestion control protocols. Related and future work is presented in Section 5, and we conclude in Section 6.

2 Multi-rate Multicast Max-Min Fairness

In this section, we present the formal network model used to examine the max-min fairness of multicast sessions, and identify a set of desirable fairness properties of max-min fair allocations for unicast networks. We then show that in this network model, the max-min fair allocation in a multicast allocation can achieve all of these desirable properties only if the sessions are multi-rate. For the reader’s convenience, a list of all the variables used is provided in Appendix A. A session, S_i , is a tuple $(X_i, \{r_{i,1}, \dots, r_{i,k_i}\})$ of session members: X_i is the session sender that transmits data within a network; each $r_{i,k}$ is a receiver that receives data from X_i . Each session contains exactly one sender and at

least one receiver. We write $r_{i,s} \in S_i$ to indicate that receiver $r_{i,s}$ is a member of session S_i .¹ We consider two types of sessions:

- If S_i is a *single-rate* session, then data must be transmitted to all receivers in S_i at the same rate.
- If S_i is a *multi-rate*, then receivers within S_i can receive data at multiple rates.

A *network graph*, G , consists of a set of nodes connected together by n links in some arbitrary fashion. The links are labeled l_1, \dots, l_n . Each link l_i has a *capacity*, c_i , that limits the aggregate rate of flow it can transmit in either direction between the two nodes it connects.² We define a network, $N = (G, \{S_1, \dots, S_m\}, \tau, \sigma)$ to be a tuple containing a network graph, G , a set of sessions, $\{S_1, \dots, S_m\}$, a mapping, τ , that maps each member of each session to a node in the network graph, and a second mapping, σ , that maps each session S_i to its type. We write $\sigma(S_i) = \mathcal{S}$ to indicate that session S_i is single-rate, and $\sigma(S_i) = \mathcal{M}$ to indicate that session S_i is multi-rate.

The mapping, τ , of a session onto the network graph has one restriction: no two members of a single session are mapped to the same node. However, there is no restriction that forbids two members of different sessions to be mapped to the same node. The network employs a routing algorithm, such that for each receiver $r_{i,k} \in S_i$, there is a sequence of links $(l_{j_1}, \dots, l_{j_s})$ that carries data from X_i to $r_{i,k}$. We refer to this set of links in this sequence as the receiver's *data-path*. The data-path for a session is defined to be the set of all links that carry data to any receiver within the session.

For a network N , we define $R_{i,j}$ to be the set of receivers in session S_i whose data-path includes link l_j , and define R_j to be the set of all receivers whose data-path includes link l_j , i.e., $R_j = \cup_i R_{i,j}$. We account for the fact that session S_i might choose a maximum rate, α_i , at which it will transmit data (α_i can be infinite). An *allocation* is an assignment of receiver rates within a network. Once an allocation has been determined, we use $a_{i,k}$ to represent the rate at which data is transmitted to receiver $r_{i,k}$ (that equals the rate at which the data is received by $r_{i,k}$, barring loss). We let $u_{i,j}$ represent an absolute measure of bandwidth (e.g., in bytes/sec) used by session S_i on link l_j to transmit data to its receivers, and u_j the amount of bandwidth used by all sessions across link j , $u_j = \sum_{i=1}^m u_{i,j}$. We refer to $u_{i,j}$ as the *session link rate* of l_j for session S_i , and u_j simply as the *link rate* of l_j . Since bandwidth for each flow is non-negative, we have $0 \leq u_{i,j} \leq u_j$. We say a link is *fully utilized* if the total bandwidth used by all sessions across the link matches its capacity, i.e., l_j is fully utilized iff $u_j = c_j$.

We require that $u_{i,j} \geq a_{i,k}$ whenever $r_{i,k} \in R_j$, i.e., any bandwidth received by a receiver must traverse its data-path. In this section, we make an additional assumption that $u_{i,j} = \max\{a_{i,k} : r_{i,k} \in R_{i,j}\}$, which is the minimum value for $u_{i,j}$ that satisfies the above requirement.³ In later sections, we examine the implications if $u_{i,j}$ is larger than this value. The assumption also allows us to model a unicast session as either a multi-rate session with a single receiver, or as a single-rate session with a single receiver. Thus, any results given in this section for networks containing a mix of single-rate and multi-rate sessions also holds for networks that contain a mix of single-rate, multi-rate, and unicast sessions.

An allocation is *feasible* if each receiver $r_{i,k}$ is assigned a rate $0 \leq a_{i,k} \leq \alpha_i$, and all receivers can receive at these rates without overutilizing any link's capacity in the network, i.e., $\forall i, k, 0 \leq a_{i,k} \leq \alpha_i$, and $\forall j, u_j \leq c_j$.⁴ The additional requirement imposed on each single-rate session S_i that all of its receivers' rates must be equal means that for any pair of receivers, $r_{i,k}, r_{i,k'} \in S_i$, when $\sigma(S_i) = \mathcal{S}$, then $a_{i,k} = a_{i,k'}$ holds. When S_i is a single-rate session, or a session of either type containing a single receiver (i.e., a unicast session), we can write the single rate at which all receivers within the session receive data simply as a_i .

Note that the feasibility of a particular allocation of receiver rates is a function of the link capacities of the network graph, G , the mapping τ , and also of the mapping σ . The dependence of an allocation's feasibility on σ is important: we will be examining how varying σ (i.e., varying sessions' types between single-rate and multi-rate) affects which allocation within a network is max-min fair.

Definition 1 (Max-min fairness) *An allocation of receiver rates is said to be **max-min fair** if it is feasible, and for any alternative feasible allocation of rates (where for each receiver $r_{i,k}$ we define $\bar{a}_{i,k}$ as an alternative feasible rate), where $\bar{a}_{i,k} > a_{i,k}$, there is some other receiver $r_{i',k'} \neq r_{i,k}$ such that $a_{i,k} \geq a_{i',k'} > \bar{a}_{i',k'}$.*

¹We assume that each receiver is a member of a single session. A receiver that is a member of two sessions can simply be viewed as two distinct receivers.

²Assigning capacity per direction is a simple extension: simply extend a bidirectional link into two unidirectional links.

³The reader that is familiar with layered approaches should see that if there is no restriction on the number of layers that a session can use, such a session link rate is easily achieved using a layered approach.

⁴Hence, in this section, we require $u_j = \sum_i u_{i,j} = \sum_i \max\{a_{i,k} : r_{i,k} \in R_{i,j}\} a_{i,k} \leq c_j$.

In other words, if any receiver $r_{i,k}$'s rate is increased beyond its max-min fair rate to obtain some other feasible allocation, then there is some other receiver whose max-min fair rate is no larger than that of $r_{i,k}$, and whose adjusted rate (to account for the increase in $r_{i,k}$'s rate) must be decreased.

When all sessions within N are single-rate (i.e., $1 \leq i \leq m, \sigma(S_i) = \mathcal{S}$), we say that N is a single-rate network, and the max-min fair allocation is called the single-rate max-min fair allocation. A similar naming convention holds when all sessions are multi-rate. The definition of max-min fairness in [17] holds only for single-rate networks,⁵ and involves a comparison of session rates rather than receiver rates as in our definition. It is easy to show that the max-min fair allocation in a single-rate network is identical under both definitions. In a network that contains multi-rate sessions, their definition is not well defined.

Just as there is always one and only one unicast max-min fair allocation [2] and one and only one single-rate max-min fair allocation [17], there is one and only one multi-rate max-min fair allocation. In fact, for any choice of σ , the network has one and only one max-min fair allocation. We show the existence of a max-min fair allocation for a network with an arbitrary σ by constructing an algorithm that achieves a max-min fair allocation for that network. The algorithm and the proof that the resulting allocation is max-min fair can be found in Appendix B; uniqueness is given by Corollary 5 in Appendix C.

Let us first examine some desirable properties of a unicast max-min fair allocation, i.e., a max-min fair allocation in a network where all sessions are unicast. It is well known that the following properties hold for a unicast max-min fair allocation [2].

Unicast Fairness Property 1 (Unicast Max-min Fairness) *For each session S_i , $1 \leq i \leq m$, either $a_i = \alpha_i$, or else there is at least one fully utilized link, l_j , where for all $1 \leq i' \leq m$, $0 < u_{i',j} \leq u_{i,j}$ (or, equivalently for the unicast case, $a_{i'} \leq a_i$ whenever $r_{i'} \in R_j$).*

Unicast Fairness Property 2 (Unicast Same Path Receiver Fairness) *If two unicast sessions, S_i and $S_{i'}$, within a unicast network have identical data-paths, then either $a_i = \alpha_i < a_{i'}$, or $a_{i'} = \alpha_{i'} < a_i$, or $a_i = a_{i'}$.*

Let us consider what makes these fairness properties desirable. To do this, we consider two *perspectives* of fairness of an allocation. From a receiver perspective, an allocation should be fair to receiver rates: a receiver's rate should be as large as possible without "stealing" bandwidth from receivers with lower rates. This is guaranteed by Unicast Property 1: there is no unused available bandwidth since some link on the receiver's data-path is fully utilized. Since there is a fully utilized link over which the receiver receives at as high a rate as any other receiver whose data-path crosses the link, increasing its rate further would result in "stealing" bandwidth from these other receivers. From a session perspective, a link's capacity should be used "fairly" by sessions. In other words, a session's allocation on a link should be as large as possible without "stealing" bandwidth from other sessions that utilize the link.

For a unicast network, the receiver and session perspectives are identical because a session's data-path is identical to its receiver's data-path, and the share of bandwidth used on each link by the session equals the receiving rate of its receiver. This is not always true in a multicast network: a receiver's data-path is only part of the session's data-path, and, in a multi-rate session, when two receivers within the session receive at different rates, there is some pair of links that have differing session link rates for that session. Hence, an allocation might be "fair" from the session perspective without being "fair" from the receiver perspective, or vice versa. One possibility is to only consider fairness properties from a single perspective (e.g., [17] considers only the session perspective). However, in this section we will assume that it is more desirable to satisfy fairness properties from both perspectives. We extend the properties of a unicast max-min fair allocation, as described in Unicast Properties 1 and 2, to multicast networks from both a session and receiver perspective.

Before presenting the desirable fairness properties for multicast networks, we introduce an example network that we will use to illustrate these different properties. Figure 1 presents a simple network with three sessions; sender X_1 in session S_1 sends to a single receiver, $r_{1,1}$, in session S_2 , sender X_2 sends to two receivers $r_{2,1}$ and $r_{2,2}$, and in session S_3 , sender X_3 sends to two receivers, $r_{3,1}$ and $r_{3,2}$. The receiving rate of a receiver, $a_{i,k}$, is indicated to the immediate right of the receiver. Each link l_j has its capacity indicated next to the link labeling, separated by a colon (e.g., $l_1 : 5$ means that $c_1 = 5$). Adjacent to the link labeling for each l_j are the session link rates, appearing in the form, $(u_{1,j} : u_{2,j} : u_{3,j})$.

⁵[17] also permits a multicast session to consist of distinct unicast connections. We model this inherently via separate unicast sessions. Such a session differs significantly from a multi-rate session achieved through layering.

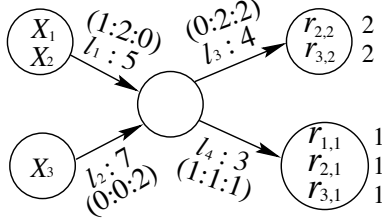


Figure 1: A sample network

Fairness Property 1 (Fully-Utilized-Receiver-Fairness) A receiver's rate $a_{i,k}$ is fully-utilized-receiver-fair if either $a_{i,k} = \alpha_i$, or there is at least one fully utilized link, l_j , where $r_{i,k} \in R_{i,j}$ and $a_{i',k'} \leq a_{i,k}$ for all receivers $a_{i',k'} \in R_j$. A session's allocation is defined to be fully-utilized-receiver-fair if the rate for each receiver in the session is fully-utilized-receiver-fair. An allocation of rates throughout the network is fully-utilized-receiver-fair if each session is fully-utilized-receiver-fair.

Fully-utilized-receiver-fairness is the multicast extension of Unicast Property 1's prevention of "stealing bandwidth" from other receivers. For instance, link l_3 is fully utilized, and lies on receiver $r_{2,2}$'s data-path. Because $r_{2,2}$ receives at a rate that is no less than any other receiver whose data-path traverses l_3 , its rate is fully-utilized-receiver-fair. Because all other receivers' rates in S_2 are fully-utilized-receiver-fair, session S_2 's allocation of rates is fully-utilized-receiver-fair. Because S_1 's and S_3 's allocations are also fully-utilized-receiver-fair, the allocation (of rates for the entire network) is fully-utilized-receiver-fair.

Fairness Property 2 (Same-Path-Receiver-Fairness) A pair of receivers $r_{i,k}$ and $r_{i',k'}$ are same-path-receiver-fair if their data-paths traverse the same set of links ($r_{i,k} \in R_j \iff r_{i',k'} \in R_j$), and either one receiver's rate is constrained by its session's maximum desired rate (i.e., either $a_{i,k} = \alpha_i < a_{i',k'}$ or $a_{i',k'} = \alpha_{i'} < a_{i,k}$), or else $a_{i,k} = a_{i',k'}$.

Same-path-receiver-fairness states that if two receivers' data-paths traverse identical links, then the receivers should receive at identical rates. In Figure 1, receivers $r_{1,1}$ and $r_{2,1}$ are a pair of receivers whose rates are same-path-receiver-fair. The reader should note that same-path-receiver-fairness is also a property of TCP-fairness [10]. If S_1 is a unicast TCP session, then, in order for $r_{2,1}$'s rate to be TCP-fair, same-path-receiver-fairness must hold for these two receivers.

Fairness Property 3 (Per-Receiver-Link-Fairness) A session S_i 's allocation is per-receiver-link-fair if for each receiver $r_{i,k} \in S_i$, either 1) $a_{i,k} = \alpha_i$, or 2) there is a link l_j that is fully utilized ($\exists j, r_{i,k} \in R_j, u_j = c_j$), and for other sessions $S_{i'}, u_{i',j} \leq u_{i,j}$. An allocation of rates throughout the network is per-receiver-link-fair if each session's allocation is per-receiver-link-fair.

Fairness Property 4 (Per-Session-Link-Fairness) An allocation is per-session-link-fair for a session S_i if $a_{i,k} = \alpha_i$ for each receiver in S_i or there exists a fully utilized link l_j in S_i 's data-path where for other sessions $S_{i'}, u_{i',j} \leq u_{i,j}$. An allocation of rates throughout the network is per-session-link-fair if each session's allocation is per-session-link-fair.

Per-receiver-link-fairness requires that session S_i gets a "fair share" of link rate along every path from sender X_i to its receivers. Per-session-link-fairness is a weaker version of this: a session must get a "fair share" of link rate on at least one link in its data-path (i.e., along the data-path of at least one receiver). In Figure 1, session S_2 is per-session-link-fair: on the data-path to receiver $r_{2,2}$, link l_3 is fully utilized and session S_2 's link rate on l_3 is no less than the link rates of other sessions on l_3 . It is also per-receiver-link-fair, because similar conditions hold on the data-path of its other receiver, $r_{2,1}$. Sessions S_1 and S_3 are also both per-receiver-link-fair and per-session-link-fair, making the network allocation both per-receiver-link-fair and per-session-link-fair.

It is fairly easy to see that in a unicast network, Fairness Property 2 and Unicast Property 2 are identical, and the remaining multicast fairness properties are identical to Unicast Property 1. We now proceed to establish properties of max-min fair allocations in terms of the types of sessions (multi-rate or single-rate) within the network. Any non-trivial proofs can be found in Appendix C.

Theorem 1 *A multi-rate max-min fair allocation satisfies the Fairness Properties 1, 2, 3, and 4. In other words, the multi-rate max-min fair allocation is fully-utilized-receiver-fair, same-path-receiver-fair, per-receiver-link-fair, and per-session-link-fair.*

Theorem 1 tells us that if all sessions are multi-rate, then the max-min fair allocation satisfies all of our desired fairness properties. We now introduce a mathematical ordering among allocations that allows us to comparatively examine the “max-min fairness” of an allocation within a network:

Definition 2 *We say a vector (x_1, x_2, \dots, x_k) is ordered if for all $i, 1 \leq i < k, x_i \leq x_{i+1}$. Let $X = (x_1, x_2, \dots, x_k)$ and $Y = (y_1, y_2, \dots, y_k)$ be ordered vectors. We write $X \dashv Y$ (and say X is min-unfavorable to Y) if no i exists such that $x_i > y_i$, or for any i where $x_i > y_i$, there is some $j < i$ where $x_j < y_j$.*

Note that under the above definition, \dashv is reflexive ($X \dashv X$), non-symmetric ($X = Y \iff X \dashv Y \wedge Y \dashv X$), and transitive ($W \dashv X \wedge X \dashv Y \implies W \dashv Y$). Furthermore, for any pair, X and Y , of ordered vectors of identical length, either $X \dashv Y$ holds, or $Y \dashv X$ holds, or both. Min-unfavorability is similar to alphabetizing two text strings of the same length. Let x_i represent the i th character of the first string, and y_i represent the i th character of the second string. Then $X \dashv Y$ if and only if $X = Y$ or an alphabetization places X before Y . A more general version of this ordering has been applied specifically within unicast networks [5]. Let us now see how this ordering relation relates to multicast max-min fairness:

Lemma 1 *Let $A = (a_1, \dots, a_s)$ be the ordered vector of receiver rates in a max-min fair allocation in a network $N = (G, \{S_1, \dots, S_m\}, \tau, \sigma)$, and let $B = (b_1, \dots, b_s)$ be the ordered vector of receiver rates for some other feasible allocation in N . Then $B \dashv A$.*

Note that the network N in Lemma 1 can have any arbitrary session type mapping, σ (i.e., some sessions can be multi-rate, while others are single-rate). However, σ must be fixed when applying the lemma. Lemma 1 along with the definition of min-unfavorability can be combined to show that the max-min fair allocation maximizes the minimum rates in a network: since all allocations are min-unfavorable to the max-min fair allocation, there exists a threshold rate x' such that for any rate $z < x'$, the number of receivers that receive at or below z is minimal (smaller or equal) within the max-min fair allocation. Furthermore, the number of receivers that receive at or below x' is minimized (strictly smaller) within the max-min fair allocation. This result can be stated more formally as a general property of min-unfavorability:

Lemma 2 $X \neq Y, X \dashv Y \iff \exists x' \text{ such that } \forall z < x', |\{x_i \in X : x_i \leq z\}| \geq |\{y_i \in Y : y_i \leq z\}| \text{ and } |\{x_i \in X : x_i \leq x'\}| > |\{y_i \in Y : y_i \leq x'\}|.$

Because the min-unfavorable relation is transitive, it gives a strict ordering among the feasible allocations for a network, where the max-min fair allocation is the maximum under the ordering. Thus, one can quantitatively compare the max-min fairness of two allocations A and B , where $A \dashv B$ means that B is “more max-min fair” than A , and the minimum receiver rates are larger in B than in A .⁶

2.1 Fairness limitations of single-rate sessions

Theorem 1 states that a multi-rate max-min fair allocation satisfies our four desirable fairness properties. Let us now see where a single-rate max-min fair allocation fails to do so. The fact that single-rate max-min allocation is per-session-link-fair is a direct consequence of the results in [17]. However, the single-rate max-min fair allocation can fail to satisfy the other fairness properties. Consider the simple example in Figure 2, whose labeling is performed in an identical manner to that of Figure 1. Here, we have a network with two sessions, S_1 and S_2 , whose respective senders, X_1 and X_2 , are located at the same point in the network. We assume that the maximum desired rates are large, $\alpha_1 = \alpha_2 = 100$, such that they do not bound receiving rates in this network. Session S_1 is a single-rate session containing three receivers $r_{1,1}, r_{1,2}, r_{1,3}$, session S_2 is a unicast session whose receiver $r_{2,1}$ is located at the same point in the network as receiver $r_{1,1}$. In the max-min fair allocation, receivers in session S_1 receive at a rate of 2 (since

⁶If one prefers to think in terms of utility rather use an ordering relation, it is fairly easy to construct a utility function, U , for allocations within a network, such that for any two allocations $A \neq B$, it is the case that $U(A) < U(B) \iff A \dashv B$. For such a utility function, the max-min fair allocation is Pareto-optimal [15].

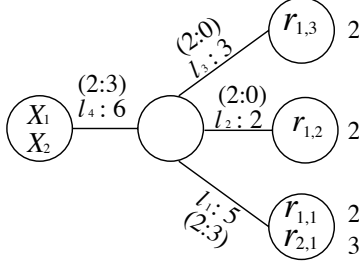


Figure 2: An example where a single-rate session would fail all but one of the fairness properties.

this fully utilizes link l_2 and all receivers must receive at the same rate in a single-rate session), the receiver in S_2 receives at a rate of 3. Receivers $r_{1,1}$ and $r_{2,1}$ fail to achieve same-path-receiver-fairness, since they have the same data-paths, but differing receiving rates. Receiver $r_{1,3}$'s rate does not satisfy fully-utilized-receiver-fairness, because there is no fully utilized link along its data-path on which its rate is the largest compared to other receivers whose data-paths cross the same link. It follows that fully-utilized-receiver-fairness does not hold for session S_1 , nor does it hold for the network. Last, per-receiver-link-fairness fails to hold for session S_1 (hence for the network as well) on the data-path to receiver $r_{1,3}$, since no link on this data-path is fully utilized. Per-receiver-link-fairness also fails to hold on the data-path to receiver $r_{1,1}$, since link l_1 is the only fully utilized link on its path and the link rate for session S_1 on that link is smaller than that of session S_2 . This example indicates that three out of the four desirable properties can fail to hold for single-rate max-min fair allocations.

We have examined the extent to which our four desirable properties hold for networks in which all sessions are the same type. Let us now consider these properties in the context of a network that contains a combination of multi-rate and single-rate sessions. Single-rate sessions are likely to always exist due to application constraints, such as a requirement that all receivers must complete receipt of data at approximately the same time.

Theorem 2 Consider a network $N = \{G, \{S_1, \dots, S_m\}, \tau, \sigma\}$ in which session types can differ, i.e., there can exist a pair of sessions, $S_i, S_{i'} \in N$ such that $\sigma(S_i) \neq \sigma(S_{i'})$. Then, the following are properties of the max-min fair allocation of N :

- (a) Fully-utilized-receiver-fairness holds for each receiver $r_{i,k} \in S_i$ where $\sigma(S_i) = \mathcal{M}$.
- (b) per-receiver-link-fairness holds for each session S_i where $\sigma(S_i) = \mathcal{M}$.
- (c) Per-session-link-fairness holds for all sessions S_i .
- (d) Same-path-receiver-fairness holds between any two receivers $r_{i,k}$ and $r_{i',k'}$ where $\sigma(S_i) = \sigma(S_{i'}) = \mathcal{M}$.
- (e) If $\sigma(S_i) = \mathcal{M}$ and $\sigma(S_{i'}) = \mathcal{S}$, and $r_{i,k} \in S_i$ and $r_{i',k'} \in S_{i'}$ have identical data-paths, then either $a_{i,k} = \alpha_i$ or $a_{i,k} \geq a_{i',k'}$.

Theorem 2 states that, even with single-rate sessions within the network, all four desirable fairness properties continue to hold for session link rates of multi-rate sessions, and for receiver rates of receivers belonging to multi-rate sessions within the max-min fair allocation. Hence, multi-rate sessions maintain their desirable fairness properties even when there are single-rate sessions within the network.

Let us also examine another way in which multi-rate multicast makes the max-min fair allocation for a network “more max-min fair”. Recall that if an allocation A is min-unfavorable to an allocation B , then B is “more max-min fair” than A . Let us now consider how the max-min fair allocations compare for any two networks, N and \bar{N} that differ only in their sessions' types.

Lemma 3 Let $N = (G, \{S_1, \dots, S_m\}, \tau, \sigma)$ and $\bar{N} = (G, \{S_1, \dots, S_m\}, \tau, \bar{\sigma})$ be networks where the set of multi-rate sessions in \bar{N} is a subset of the set of multi-rate sessions in N , (i.e., $\forall i, \bar{\sigma}(S_i) = \mathcal{M} \Rightarrow \sigma(S_i) = \mathcal{M}$). If A is the ordered vector of max-min fair receiver rates in N , and \bar{A} is the ordered vector of max-min fair receiver rates in \bar{N} , then $\bar{A} \dashv A$.

Proof: Since an allocation for a single-rate session is feasible for a multi-rate session, \bar{A} is a feasible allocation in N . Since A is the max-min fair allocation in N , by Lemma 1, $\bar{A} \dashv A$. ■

Lemma 3 tells us that as we “replace” single-rate sessions with identical multi-rate sessions (i.e., the only difference between the single-rate session and its replacement is the session type), then the max-min fair allocation is “more max-min fair”. Hence, the “most max-min fair” allocation is the one in which all sessions are multi-rate:

Corollary 1 *Let $N = (G, \{S_1, \dots, S_m\}, \tau, \sigma)$ be a multi-rate network ($\forall i, \sigma(S_i) = \mathcal{M}$), and let $\bar{N} = (G, \{S_1, \dots, S_m\}, \tau, \bar{\sigma})$ be the single-rate version of N , ($\forall i, \bar{\sigma}(S_i) = \mathcal{S}$). Let A be the ordered vector of receiver rates for a multi-rate max-min fair allocation within N , and let B be the ordered vector of receiver rates in \bar{N} . Then $B \preceq A$, and if $A \neq B$, then $A \not\preceq B$.*

Last, let us consider how varying session types affects receiving rates on a session-by-session basis. We can prove that if all sessions’ types are fixed except for session S_i , then if S_i is multi-rate, all of its receivers will receive at rates that are no less than what they would receive at if S_i is single-rate (see Lemma 9 in Appendix C). Unfortunately, this result does not extend to the case when several sessions can switch types. In fact, it is rather difficult to say what happens to receiver rates due to changes in the session type or the network topology. For example, one might conjecture that removing a receiver $r_{i,k}$ from a session would only increase other receiver fair rates. Our intuition was that since the removal frees up bandwidth that can then be used by other receivers whose data-path crosses $r_{i,k}$ ’s data-path. However, the reallocation of bandwidth after the receiver is removed can cause receiver rates (both in session S_i and in other sessions) to vary in either direction.

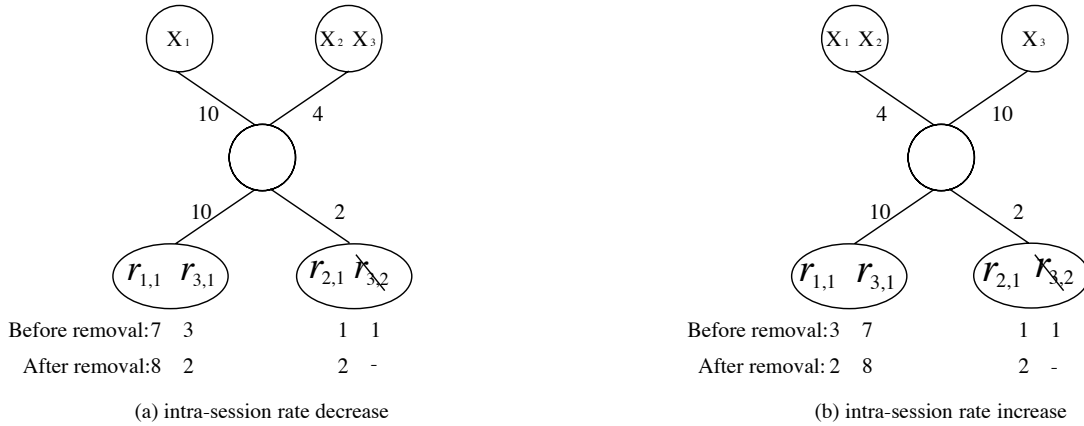


Figure 3: The change in max-min fair rates due to a removal of a receiver from a session.

To see this, consider the examples in Figure 3. Both networks contain three multi-rate sessions, S_1 , S_2 , and S_3 . S_1 and S_2 each contain a single receiver, S_3 contains two receivers, the second ($r_{3,2}$) is subsequently removed. The max-min fair rates for receivers are indicated before and after this removal. Note that in Figure 3(a), $r_{3,1}$ ’s max-min fair rate decreases and $r_{1,1}$ ’s rate increases as a result of the removal. In Figure 3(b), $r_{3,1}$ ’s rate increases and $r_{1,1}$ ’s rate decreases. These figures demonstrate that removing receivers from sessions can have a nonobvious impact on the max-min fair rates of the remaining receivers in the network. Additional results appear in Appendix D.

We now summarize the main results of this section. We have shown that if multicast sessions are multi-rate, then the max-min fair allocation is “more max-min fair” than if the sessions are restricted to being single-rate. We have demonstrated this by showing that there are four desirable fairness properties that hold in the multi-rate max-min fair allocation that do not necessarily hold in a single-rate max-min fair allocation. We also examined networks in which some of the sessions are single-rate, while the remaining are multi-rate. By examining fairness properties on a per-session basis, we find that all of the fairness properties hold in general only in multi-rate sessions. Last, we use the min-unfavorable relation to comparatively examine which of any two allocations for a network is “more max-min fair”. We find that “replacing” single-rate sessions by multi-rate sessions makes the max-min fair allocation more “max-min fair”, which means that when all sessions are multi-rate, the max-min fair allocation is the “most max-min fair”.

3 Achieving Multi-Rate Max-Min Fairness with Layering

In the previous section, we motivated the use of multi-rate sessions by showing that they yield more desirable max-min fair allocations. One way to then obtain these rates in practice is to have the sender configure layers so that each receiver can obtain its fair rate by joining some subset of layers. However, the number of layers can be as large as the number of receivers in the session, making such an approach infeasible for large multicast sessions. Furthermore, the number of layers and the rate per layer is often beyond the control of the session itself, due to application-specific requirements, a limitation in the availability of multicast groups, or because it is too difficult for the sender to obtain the feedback needed to appropriately configure the rates of each of the layers. In this section, we examine how receivers can obtain their long term average max-min fair rates by repeated joins and leaves from multicast groups on which data is sent at a restricted set of rates. We will see that such a mechanism will force us to reconsider our previous assumption of how receiver rates impact link rates in the network.

Let us first discuss the implementation of a layered multicast approach. Data to be transferred is split into M layers by the sender, where layers are transmitted on separate multicast groups, each at some rate. The layers are ordered L_1, \dots, L_M , such that all receivers desiring transmission join the group containing layer L_1 , and any receiver that joins the group containing layer L_j must also join or already be joined to layer L_i for all $1 \leq i < j$ (henceforth, this is implied when we say that the receiver joins the layer or joins *up to* the layer). A receiver joined up to layer L_i receives data from the sender at an aggregate rate equal to the sum of the rates of layers L_1 through L_i . Joining layers increases the aggregate rate, while leaving layers decreases the aggregate rate.⁷

Let us examine why receivers must join and leave layers to obtain their fair rates. An obvious alternative is to require receivers to choose rates that can be obtained by joining up to a given layer and remaining at that rate for the duration of the session. This makes a finite set of rates available to the receiver. However, if these layers cannot be configured to the needs of receivers for reasons described above, the max-min fair allocation might not even exist! As an example, consider a simple network that consists of a single link with capacity c , and let there be two layered multicast sessions, S_1 , and S_2 that traverse this link. Each session contains a single receiver, respectively denoted r_1 and r_2 . The sender for session S_1 provides three layers, and sends at a rate of $c/3$ per layer. The sender for session S_2 provides two layers, and sends at rate $c/2$ per layer. The set of feasible allocations is $\{(0, 0), (0, c/2), (0, c), (c/3, 0), (c/3, c/2), (2c/3, 0), (c, 0)\}$, where (a_1, a_2) implies receiver r_i receives at a rate of a_i . None of these allocations are max-min fair. For instance, $(a_1, a_2) = (c/3, c/2)$ is not max-min fair since $(\bar{a}_1, \bar{a}_2) = (2c/3, 0)$ is feasible, and $a_1 < \bar{a}_1$, but $a_2 > \bar{a}_2$, hence there is no j where $\bar{a}_j < a_j \leq a_1$ ⁸ (contradicting the defined requirement for max-min fairness). The reader can easily verify that none of the other feasible allocations is max-min fair.

Although it is not possible to achieve a max-min fair rate allocation when receivers are restricted to joining a fixed set of layers for the entire length of a session, it is possible to achieve long-term average max-min fair rates through joins and leaves. The idea of using long term average rates also appears in current definitions of TCP-fairness [10, 18, 3, 13]. We define the *quantum*, Δt , to be the minimum amount of time over which a receiver's average rate is computed. We say that a rate of r is obtained through a link during the i th quantum if $r\Delta t$ bytes pass through the link between times $i\Delta t$ and $(i+1)\Delta t$. We say that a link l_j can support a capacity of c_j if it is able to forward $c_j\Delta t$ bytes within each time quantum.

Let us now consider an idealized network where a receiver can use joins and leaves to obtain its fair rate. The network is ideal in that we assume that network propagation delays and leave latencies are negligible compared to Δt and to packet inter-arrival times for each session. In this model, a packet traverses a link l_j only if it is received by some receiver $r_{i,k} \in R_j$. We also assume that all packets are of equal size, and for any receiver $r_{i,k}$, let $a_{i,k} \leq \alpha_i$ be its fair packet rate (in packets/sec) within the network. Consider a single layer (multicast group), where the transmission rate on the layer, $\rho \geq \max\{a_{i,k} : r_{i,k} \in S_i\}$. Receiver $r_{i,k}$ joins the single layer so that it receives the first $a_{i,k}\Delta t$ packets within the quantum,⁹ then leaves the group. This is clearly possible, since $a_{i,k} \leq \alpha_i \leq \rho$, and $\rho\Delta t$ packets are transmitted on the layer during the quantum.

In this scenario, for any link l_j and session S_i where $|R_{i,j}| > 0$, there is some receiver $r_{i,k'}$ that receives $a_{i,k'} = \max\{a_{i,k} | r_{i,k} \in R_{i,j}\}$ packets per time quantum. Hence, this is the minimum number of packets that traverse link l_j

⁷We make the assumption that there is some utility in receiving at a faster rate, e.g., audio and video transmissions increase in clarity, reliable data transmissions take less time.

⁸Or less formally, r_1 's increase in rate does not result in a decrease in any receiver's rate whose original rate was less than r_1 's.

⁹If $a_{i,k}\Delta t$ is not an integer, then it can elect to receive $\lfloor a_{i,k}\Delta t \rfloor$ packets in each quantum, and periodically receive $\lceil a_{i,k}\Delta t \rceil$ to come arbitrarily close to $a_{i,k}\Delta t$.

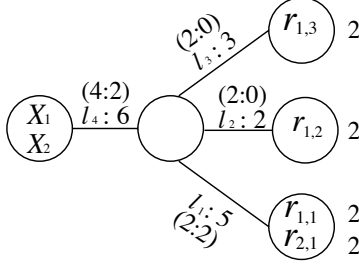


Figure 4: An example where a network fails to achieve session-perspective fairness properties due to redundancy.

for session S_i per quantum. Transmitting exactly this number of packets requires that all other receivers $r_{i,k} \in R_{i,j}$ receive a subset of the packets that are received by $r_{i,k'}$ per quantum. When this is not the case, $u_{i,j} > a_{i,k'}$.

Definition 3 We define the **redundancy** of a link l_j for a session S_i to be $u_{i,j} / \max\{a_{i,k} | r_{i,k} \in R_{i,j}\}$, where $u_{i,j}$ is the long-term average link rate l_j by session S_i , and $a_{i,k}$ is the long-term average rate for receiver $r_{i,k}$. We say a session's bandwidth utilization of a link is **efficient** for session S_i if the link's redundancy for that session is one, and define a session S_i 's **efficient link rate** to equal $\max\{a_{i,k} | r_{i,k} \in R_{i,j}\}$.

Note that our assumption in Section 2 that $u_{i,j} = \max\{a_{i,k} : r_{i,k} \in R_{i,j}\}$ amounts to an assumption that multi-rate sessions are efficient (i.e., on all links in the network, a multi-rate session's link rate equals its efficient link rate). When there are multi-rate sessions that are not efficient, a multi-rate max-min fair allocation might not satisfy per-session-link-fairness (and hence might not satisfy per-receiver-link-fairness). To show this, we consider the network shown in Figure 4, whose labeling is similar to that of Figures 1 and 2. We again assume that the maximum desired rates are large so as not to bound receiving rates, e.g., let $\alpha_1 = \alpha_2 = 100$. Here, session S_1 is multi-rate with a redundancy of 2 over the shared link, l_4 . Since the maximum receiving rate for receivers in S_1 (all of whose data-paths traverse l_4) is 2, $u_{1,4} = 4$. Since this is the only link that is fully utilized, and $u_{1,4} > u_{2,4}$, per-session-link-fairness fails. It follows that per-receiver-link-fairness fails to hold for session S_2 as well.

It is trivial to show that the fairness properties that do not compare session link rates, (specifically same-path-receiver-fairness and fully-utilized-receiver-fairness), continue to hold even when sessions are not efficient.

To understand how important coordination between receiver joins and leaves is for redundancy, let us examine what happens on a shared link when there is no implicit join/leave coordination. Assume each receiver $r_{i,k}$ within session S_i randomly chooses the $a_{i,k}\Delta t$ packets it should receive within the quantum, with each packet having an equally likely chance of being chosen as any other in that quantum. In this case, $E[U_{i,j}] = \rho(1 - \prod_{t=1}^s (1 - a_{i,k_t}/\rho))$, where $\{a_{i,k_1}, \dots, a_{i,k_s}\}$ are the rates of receivers that are members of the set $R_{i,j}$ (derivation in Appendix E).

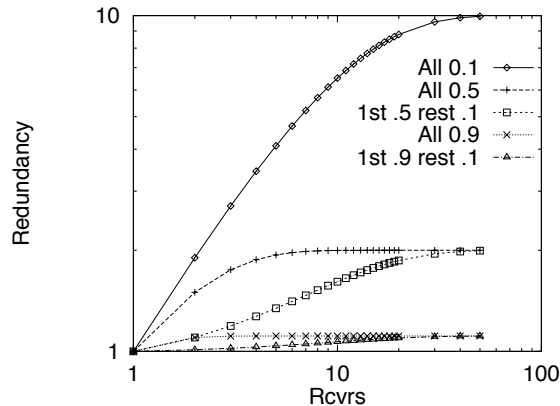


Figure 5: Redundancy of a single layer with random joins

Figure 5 shows how the number of receivers within a session that utilize a link (i.e., $|R_{i,j}|$) impacts the redundancy of a layer when each receiver randomly chooses the $a_{i,k}\Delta t$ packets it receives within the quantum. The number of

receivers is shown on the x -axis, while the session's redundancy is indicated on the y -axis. The curves represent various configurations of $\{a_{i,1}, \dots, a_{i,s}\}$. For curves labeled **All** z , ($z = 0.1, 0.5$, or 0.9), each receiver's $a_{i,k}$ is respectively set to $.1\rho$, $.5\rho$, and $.9\rho$ for each receiver. For curves labeled **1st w rest** z , $a_{i,1} = w\rho$, and $a_{i,s} = z\rho$ for $1 < s \leq |R_{i,j}|$. Note that in each plot, the efficient link rate remains constant as the number of receivers is varied.

Note that for redundancy to be high, the ratio of efficient link rate to the transmission rate (i.e., $\max_{r_{i,k} \in R_{i,j}} \{a_{i,k}\} / \rho$) must be small. In fact, the redundancy can only be as large as the multiplicative inverse of this value (e.g., $\max\{a_{i,k} | r_{i,k} \in R_{i,j}\} / \rho = .1$ yields a redundancy of 10), and asymptotically reaches this value with an increase in the number of receivers that share the link. In other words, for redundancy to be high, all receivers must require only a small percentage of packets from a layer.

A second result is that for a fixed efficient link rate, redundancy increases most rapidly as a function of the number of receivers when all receivers receive at the same rate of a_{\max} , equal to the efficient link rate. In other words, an upper bound on how additional receivers impact redundancy is obtained by considering a network in which all receivers within a session have identical fair rates.

These results give a preliminary indication as to what impacts the magnitude of redundancy within a network. We find that having additional layers often leads to a reduction in redundancy that is sometimes substantial, and that it never increases redundancy beyond that exhibited for the single-layer case. Details of these results can be found in Appendix E.

3.1 The impact of redundancy on fair rates

Let us now examine the impact that redundancy has on fairness within a network. We now demonstrate why sessions with lower redundancy are “more max-min fair” than corresponding ones with high redundancy. We begin by relaxing our assumption made in Section 2 that $u_{i,j} = \max\{a_{i,k} : r_{i,k} \in R_{i,j}\}$. We extend our definition of a session to be a tuple $S_i = (X_i, \{r_{i,1}, \dots, r_{i,k_i}\}, v_i)$ that now includes a redundancy function v_i . Here, v_i maps a set (of arbitrary size) of receiver rates to a link rate. Given an allocation of receiver rates, A , session S_i 's link rate for link l_j is computed as $u_{i,j} = v_i(\{a_{i,k} : r_{i,k} \in R_{i,j}\})$. In Section 2, v_i is simply the max operation. Since $u_{i,j} \geq a_{i,k}$ must hold whenever $r_{i,k} \in R_j$ (for reasons discussed in Section 2), it follows that $v_i(\{a_{i,k} : r_{i,k} \in R_{i,j}\}) \geq \max\{a_{i,k} : r_{i,k} \in R_{i,j}\}$.¹⁰

Lemma 4 *Let $N = (G, \{S_1, \dots, S_m\}, \tau, \sigma)$, and $\bar{N} = (G, \{\bar{S}_1, \dots, \bar{S}_m\}, \tau, \sigma)$ be identical networks, where each session S_i in N is identical to \bar{S}_i in \bar{N} , except for their respective redundancy functions, v_i and \bar{v}_i . Assume sessions in \bar{N} exhibit higher redundancy than those in N , (i.e., for each session S_i and any set of real numbers, X , $v_i(X) \leq \bar{v}_i(X)$). Let A be the max-min fair allocation in N and \bar{A} the max-min fair allocation in \bar{N} . Then $\bar{A} \dashv A$.*

Proof: First consider allocation \bar{A} in \bar{N} , and let $\bar{a}_{i,k}$ represent receiver $r_{i,k}$'s rate under allocation \bar{A} . Since \bar{A} is max-min fair, it is feasible in \bar{N} . Feasibility implies that for any link l_j in \bar{N} , $c_j \geq \sum_i \left(\sum_{X=\{\bar{a}_{i,k} : r_{i,k} \in R_{i,j}\}} \bar{v}_i(X) \right)$. Our assumption for that $v_i(X) \leq \bar{v}_i(X)$ for any set of real numbers X gives us that $\sum_i \left(\sum_{X=\{\bar{a}_{i,k} : r_{i,k} \in R_{i,j}\}} \bar{v}_i(X) \right) \geq \sum_i \left(\sum_{X=\{a_{i,k} : r_{i,k} \in R_{i,j}\}} v_i(X) \right)$, which is the link rate for link l_j for allocation A in network N . Since l_j has identical capacity in N and \bar{N} , it must be that A is feasible in N , and by Lemma 1, we have $\bar{A} \dashv A$. ■

Lemma 4 states the following: assume that sessions are “replaced” by sessions that are identical, except that the session link rates required to support a given set of receiver rates are higher (e.g., the amount of coordination of joins and leaves between receivers within a session is reduced). It follows that the resulting max-min fair allocation is “less max-min fair” than the max-min fair allocation for the network with the sessions prior to the “replacement”.

We know that a redundancy greater than one produces max-min fair rate allocations within the network that might not exhibit the session-perspective fairness properties, per-receiver-link-fairness and per-session-link-fairness. Also, using the min-unfavorable relation, we have shown that increased redundancy might reduce the “max-min fairness” of a max-min fair allocation. Let us now quantitatively examine how redundancy may impact fair rates. Consider a set of n sessions whose receiver rates are constrained by the same link, l with capacity c . Let m of these sessions be multi-rate with a redundancy of v on link l , and the remaining $n - m$ sessions have redundancy 1. Since we assume that all receivers' rates are constrained by link l , their max-min fair rates are all equal to $\frac{c}{(n-m)+mv}$. Figure 6 shows

¹⁰Actually, v_i must also be non-decreasing and continuous in the following sense: let A and A' be sets of rates where $|A| = |A'|$, and let $\psi : A \rightarrow A'$ be a bijection. Then if $\forall a_{i,k} \in A, \psi(a_{i,k}) \geq a_{i,k}$, then $v_i(A) \leq v_i(A')$ (non-decreasing). Also, $\forall \epsilon > 0, \exists \delta$ such that if $\forall a_{i,k} \in A, |\psi(a_{i,k}) - a_{i,k}| < \delta \Rightarrow |v_i(A) - v_i(A')| < \epsilon$ (continuity).

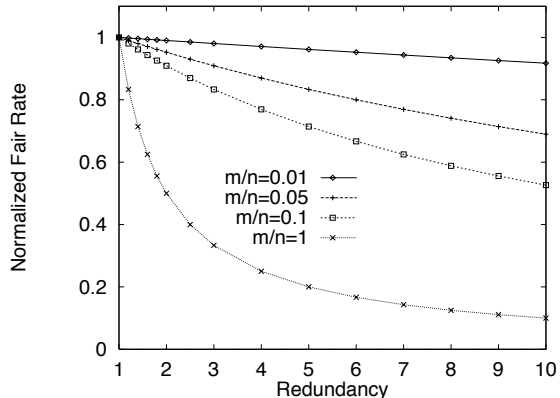


Figure 6: The impact of redundancy on fair rates.

the receivers’ rates as a function of the redundancy, v . The x -axis indicates v , the various curves represent various values of the ratio of sessions, m/n , that exhibit redundancy v . The y -axis presents the fair rate normalized by c/n , the fair rate for all the receivers in the network when all sessions are efficient.

Figure 6 indicates that even modest levels of redundancy can substantially reduce the fair receiver rates for all sessions in the network. From this we can draw two conclusions: first, it is important to maintain low redundancy on network links to keep fair rates high. Second, when multi-rate sessions make up a small percentage of sessions in the network, they have less of an impact on the fair rates of sessions. Due to the current proliferation of unicast traffic within the network, we expect that less than 5% of sessions within the network will be multi-rate. This means that low levels of redundancy greater than one can be tolerated.

These results raise an interesting dilemma: should multi-rate protocols be used to achieve fairness from the receivers’ perspectives, even if it means failing to achieve per-session-link-fairness (a fairness property that holds when all network sessions are single-rate and unicast)? We argue that yes, multi-rate protocols should still be used, because the “unfair” additional usage of link bandwidth due to redundancy can be justified in that the session is transmitting data to multiple receivers. A similar argument is used in [7] to allocate link bandwidth to sessions in a manner that is proportional to the number of receivers within the session.

The reduction in rate due to redundancy can occur whenever a multi-rate session tries to achieve some form of fairness using joins and leaves of layers. For example, in [18], receiver join experiments are coordinated within a network where TCP-fairness is the fairness criteria. The coordination prevents “bottleneck bandwidth allocated to [the] protocol instance [from] not being fully exploited.” This lack of “exploitation” is, in effect, an artifact of redundancy.

4 Redundancy in Practical Congestion Control Protocols

In Section 3, we showed that a lack of join and leave coordination within a session increases the session’s redundancy on links shared by that session’s receivers, which is likely to reduce their fair receiving rates. Our final contribution is to show that redundancy can easily be kept quite low in practice. We show this by measuring the redundancy of several Internet layered congestion control protocols that vary in the degree to which joins are coordinated among receivers. In these protocols, receivers react to congestion by leaving layers, and probe for available bandwidth by joining layers. We compute each protocol’s redundancy using analysis and simulation for simple network models. Because of the simplicity of the models, there may be some differences between what we observe and what will actually occur in practice. However, we do not expect these differences to alter results significantly enough to change our conclusion.

In each protocol, a receiver leaves the highest layer joined (unless only joined to one layer) whenever it observes a *congestion event*: an indication that some part of its data-path is being overutilized. In practice, a congestion event may be the loss of a packet by the receiver, or a bit set within a packet by the network used to indicate that the receiving rate should be lowered [12]. If no congestion events are observed by a receiver within a sequence of packet arrivals, it joins an additional layer (unless already joined to all layers). Using these protocols, a receiver repeatedly adjusts the set of layers to which it is joined for the duration of the session. The protocols differ in the degree to which joins are

coordinated within a session.

- In the *Uncoordinated* protocol, there is no inherent coordination: upon receiving a packet, a receiver randomly decides whether to join an additional layer.
- In the *Deterministic* protocol, there is also no inherent coordination; a receiver joins an additional layer after receiving a fixed number of packets without loss since its last join or leave event.
- In the *Coordinated* protocol, the sender indicates (e.g. through a field within its transmitted packet) when receivers should join an additional layer. This is done in such a way so that when the field indicates that receivers joined up to layer i should join layer $i + 1$, it also indicates that receivers joined up to layer $j < i$ should join layer $j + 1$.

The additional details of the protocols (layer rates, join-period) are based on the choices made in [18]. For instance, we require that the aggregate rate of layers 1 through i equals 2^{i-1} , and that the expected number of packets received by a receiver between a previous join/leave event to its join to layer $i + 1$ equals $2^{2(i-1)}$.¹¹ Because of these protocols’ similarities to the protocol in [18], we anticipate these protocols are suitable for the same set of continuous stream and reliable bulk data transfer applications described in [18]. Due to a lack of round-trip-time dependence, these protocols come closer to achieving max-min fair rates than TCP-fair rates. See Appendix F for a more precise description of these protocols and how they differ from the protocol in [18].

We model packet loss (or equivalently, congestion marking of packets as a Bernoulli loss process. The reader can consider the loss process to be fairly accurate for a network where the number of flows across links is large, so that there is little correlation between the rate of an individual flow and the link loss rate [20]. Our model also assumes that receivers’ reactions to coordinated events (shared loss, coordinated joins) take effect at the same time: two receivers that see identical loss patterns would be joined to the same set of layers. Under these conditions, it can be argued that these protocols come “close” to achieving the max-min fair rates, i.e., the expected rate does not exactly equal the max-min fair rate, but the difference is fairly small.



Figure 7: Network models for coordination experiments

Our experiments use modified star networks, as shown Figure 7, to examine how shared loss (i.e., loss on the shared link abutting the sender) and independent loss (i.e., loss on the fanout links) impact redundancy. The initial set of experiments uses the topology in Figure 7(a). Using Markov models of the protocols over this network, we examine how different values of shared and independent loss impact the redundancy of a session on the shared link. The details of these models appear in Appendix F We summarize the most important finding: redundancy is highest when receivers experience the same end-to-end loss rates. This result follows intuitively from our observation in Section 3 that redundancy is highest when all receivers receiving rates are equal to a_{\max} .

Our Markov models are too computation-intensive to allow us to examine sessions with large sets of receivers. Instead, we turn to simulation. Figure 8 shows simulations of the protocols using 8 layers with 100 receivers in the session with identical end-to-end loss rates, configured in the modified-star topology of Figure 7(b). In Figure 8(a), the shared loss rate is fixed to 0.0001 (i.e., very low shared loss), and the loss rate on each of the fanout links is given on the x -axis. Each curve shows the redundancy for the three protocols we consider.¹² Figure 8(b) plots similar results, but where the shared loss rate is .05. We see that for all protocols, redundancy remains fairly low (below 5)

¹¹In [18], the number of packets received equals $2^{2(i-1)}$ (i.e., it is a deterministic value).

¹²Each point plotted is the mean of 30 experiments where the sender transmits 100,000 packets, the variance is less than 1% with 95% confidence.

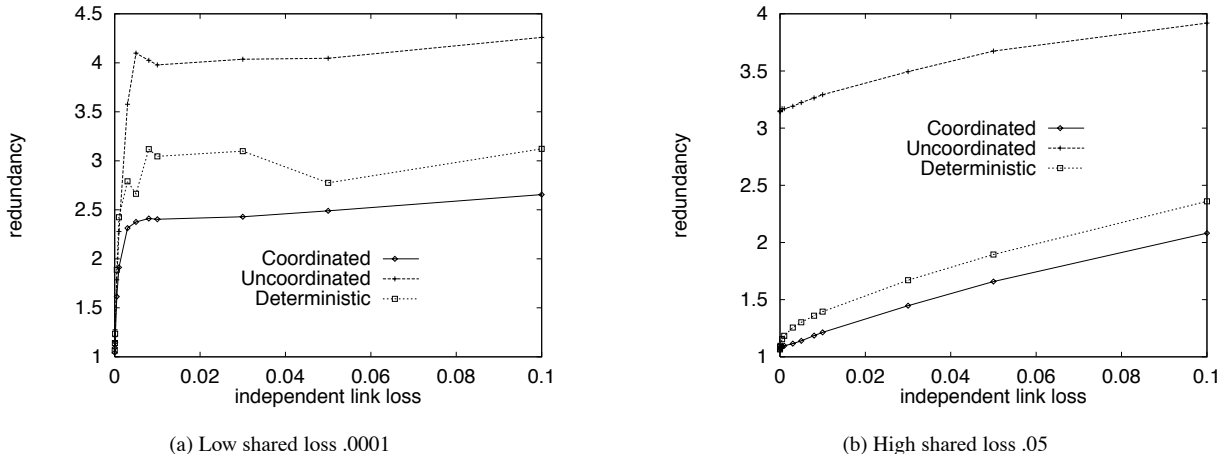


Figure 8: 100 receivers, 8 layers

for reasonable loss rates. By having the sender coordinate joins as in the Coordinated protocol, redundancy remains below 2.5, even when there are 100 receivers within the session, each of whose data-path contains the shared link. We observed negligible changes in the results when we increased the number of receivers beyond 100. Since our previous results indicate that redundancy is highest when all receivers have identical end-to-end loss rates, we can conclude that sender-coordinated congestion control protocols can keep redundancy below 2.5. This is low enough so that, in networks where multi-rate protocols make up a small percentage of sessions, multi-rate protocols will yield fair allocations with sufficiently desirable fairness properties.

5 Related / Future Work

The application of layering, in the context of video transmission, to maximize usage of available bandwidth and the benefits of coordination of receiver join events within a session are discussed in [11], and further explored in [8]. Clever use of parity coding techniques extend layering’s applicability to reliable multicast [4, 18, 14]. Preliminary experiments and definitions of various forms of fairness for layered approaches are explored in [8, 18], as well as in [9], which discusses at a high level how using a layered approach can change the max-min fair allocation. An examination that uses fairness metrics to compare various allocation strategies for layered multicast protocols is presented in [7]. There, the authors argue that link bandwidth should be allocated to sessions in some manner that is proportional to the number of receivers in the session because doing so increases the average “receiver satisfaction”. However, none of these works look consider how layered approaches affect fairness properties (in comparison to single-rate approaches) throughout a large-scale network.

Much of the remaining work that deals with multicast fairness assumes that sessions are single-rate [17, 3, 19, 13], and therefore compromise fairness from the receiver perspective, due to tight binding of receiver rates within a session. There has been some work that discusses how one might choose a single-rate session’s rate in order to maximize a measure of fairness on a per-receiver basis [6].

There are numerous issues that remain open with regard to using layering to achieve multi-rate max-min fairness. The effects of layering on desirable fairness properties for other definitions of fairness is one possible avenue for examination. We believe that many of our results can be directly applied to TCP-fairness by constructing a definition of max-min fairness where receiver rates are assigned weights (i.e., a receiver’s rate is weighted by the inverse of round trip time). It would also be interesting and useful to extend definitions of fairness to multicast sessions with multiple senders. There are also many issues that deal with the practicality of using layering to achieve fairness. One question that comes to mind is whether priority dropping schemes for layered approaches [1] might aid in reducing redundancy by increasing coordination among receivers. Also, multicast routing technology must be improved to make layered approaches practical for congestion control and fairness purposes. For instance, join and leave latencies complicate coordination among various receivers within a session, which is likely to increase redundancy. We believe

that long leave latencies will also increase redundancy (a link continues to receive at the rate prior to the leave, until the leave takes effect, while the receiver's rate reduces immediately). We expect that many such problems are solvable, perhaps with the aid of active routing technology [16]. For instance, placing the decision to add and drop layers at the active nodes, rather than at receivers, should increase the coordination of the joins and leaves of layers by downstream receivers, thereby reducing redundancy. Such an approach would make a redundancy of one feasible for a layered multi-rate session.

It is also unclear whether bandwidth can be shared fairly by sessions that measure fairness on different timescales (i.e., use different quanta), especially in networks like the Internet where a session's fair allocation may vary due to startup and/or termination of other sessions within the network. Finally, our models contain numerous simplifications of what exists in practice; they are merely used to illustrate concepts, identify challenges, and provide a basic understanding of what can be expected in practice. Extensive development and testing is still necessary to verify that our hypotheses presented here do in fact occur in practice.

6 Conclusion

We have explored how multi-rate multicast, achievable using layered multicast approaches, can impact fairness within a network. In particular, we showed that in theory, multi-rate sessions can achieve several desirable fairness properties that cannot be achieved in general networks using single-rate sessions. In a practical environment, we demonstrate how receivers can join and leave layers so that their rates are max-min fair over a long term average. Unfortunately, this join-leave process has several practical difficulties. One difficulty that we address is redundancy: an excessive use of bandwidth by a session over a link shared by multiple receivers in the session. High redundancy not only leads to failure of several fairness properties from a session perspective (i.e., fairness of session link rates), but is also likely to reduce most receivers' fair rates. Our subsequent analysis shows, however, that based on the portion of network sessions that are expected to be multi-rate, practical solutions can keep the amount of redundancy low enough such that layering can be used to improve fairness within multicast networks.

References

- [1] S. Bajaj, L. Breslau, and S. Shenker, *Uniform versus Priority Dropping for Layered Video*, Proceedings of ACM SIGCOMM'98, Vancouver, CA, September 1998.
- [2] D. Bertsekas and R. Gallager, *Data Networks*, Englewood Cliffs, NJ, Prentice-Hall, 1992.
- [3] S. Bhattacharyya, D. Towsley, and J. Kurose, *The Loss Path Multiplicity Problem for Multicast Congestion Control*, To appear in Proceedings of IEEE INFOCOM 99, New York, NY, March 1999.
- [4] J. Byers, M. Luby, M. Mitzenmacher, and A. Rege, *A Digital Fountain Approach to Reliable Distribution of Bulk Data*, Proceedings of ACM SIGCOMM'98, Vancouver, CA, September 1998.
- [5] Z. Cao and E. Zegura, *Utility Max-Min: An Application-Oriented Bandwidth Allocation Scheme*, To appear in Proceedings of IEEE INFOCOM 99, New York, NY, March, 1999.
- [6] T. Jiang, M. Ammar, and E. Zegura, *Inter-Receiver Fairness: A Novel Performance Measure for Multicast ABR Sessions*, Proceedings of ACM SIGMETRICS 98, Madison, Wisconsin, June 1998.
- [7] A. Legout, J. Nonnenmacher, and E. Biersack, *Bandwidth Allocation Policies for Unicast and Multicast Streams*, To appear in Proceedings of IEEE INFOCOM 99, New York, NY, March, 1999.
- [8] X. Li, S. Paul, and M. Ammar, *Layered Video Multicast with Retransmissions (LVMR): Evaluation of Hierarchical Rate Control*, Proceedings of INFOCOM 98, March 1998, San Francisco, CA.
- [9] X. Li, S. Paul, and M. Ammar, *Multi-Session Rate Control for Layered Video Multicast*, to appear in Proceedings of Symposium on Multimedia Computing and Networking, San Jose, CA, January 1999.
- [10] J. Mahdavi and S. Floyd, *TCP-Friendly Unicast Rate-Based Flow Control*, Note sent to e2e mailing list, January, 1997.
- [11] S. McCanne, V. Jacobson, and M. Vertterli, *Receiver Driven Layered Multicast*, Proceedings of ACM SIGCOMM 96, Stanford, CA, August, 1996.
- [12] Ramakrishnan, K.K., and Floyd, S., *A Proposal to add Explicit Congestion Notification (ECN) to IP*. RFC 2481, January 1999.
- [13] I. Rhee, N. Balaguru, G. Rouskas, *MTCP: Scalable TCP-like Congestion Control for Reliable Multicast*, To appear in Proceedings of IEEE INFOCOM 99, New York, NY, March 1999.
- [14] Luigi Rizzo and Lorenzo Vicisano, *RMDP: An FEC-based Reliable Multicast Protocol for Wireless Environments*, Mobile Computing and Communications Review, Volume 2, Number 2, April 1998.

- [15] S. Shenker, *Making Greed Work in Networks: A Game-Theoretic Analysis of Switch Service Disciplines*, Proceedings of ACM SIGCOMM'94, London, UK, August 1994.
- [16] Tennenhouse, D., Smith, J., Sincoskie, D., Wetherall, D., and Minden, G., *A Survey of Active Network Research*, IEEE Communications Magazine, January 1997.
- [17] H. Tzeng and K. Siu, *On Max-Min Fair Congestion Control for Multicast ABR Service in ATM*, IEEE JSAC, Vol. 15, No. 3, April 1997.
- [18] L. Vicisano, J. Crowcroft, and L. Rizzo, *TCP-like Congestion Control for Layered Multicast Data Transfer*, Proceedings of IEEE INFOCOM 98, San Francisco, CA, March, 1998.
- [19] H. Wang and M. Schwartz, *Achieving bounded fairness for multicast and TCP traffic in the Internet*, Proceedings of ACM SIGCOMM'98, Vancouver, CA, September 1998.
- [20] M. Yajnik, S.B. Moon, J. Kurose, and D. Towsley, *Measurement and Modeling of the Temporal Dependence in Packet Loss*, To appear in Proceedings of IEEE INFOCOM 99, New York, NY, March, 1999.

Appendix

A Summary of variables

| | |
|---|--|
| G | A network graph with n links. |
| $l_j, 1 \leq j \leq n$ | The j th link of N |
| $c_j, 1 \leq j \leq n$ | The capacity of link l_j |
| $S_i, 1 \leq i \leq m$ | The i th session in N |
| σ | a mapping onto each session S_i that indicates the session's type (\mathcal{M} = multi-rate or \mathcal{S} = single-rate) |
| $r_{i,k}, 1 \leq i \leq m, 1 \leq k \leq k_i$ | The k th of k_i receivers in session S_i |
| X_i | the single sender for session S_i |
| τ | A topology mapping that maps session members onto network nodes. |
| $N = (G, \{S_1, \dots, S_m\}, \tau, \sigma)$ | a network |
| α_i | The maximum desired rate for session $S_i, 0 < \alpha_i \leq \infty$ |
| $R_{i,j}$ | The set of receivers in S_i whose data-path traverses l_j |
| R_j | The set of receivers over all sessions whose data traverses l_j |
| $a_{i,k}$ | The data rate for transmission to receiver $r_{i,k}$ |
| r_i | The receiver in a unicast session S_i |
| a_i | The data rate in a unicast or single-rate session |
| $u_{i,j}$ | The link rate for session S_i on link l_j |
| u_j | The link rate for link l_j (i.e., $\sum_i u_{i,j}$) |
| Defined in Section 3: | |
| ρ | The aggregate rate of the "single-layer" |
| v_i | A more general session link rate function |

B Max-min fair construction and existence proof

The following algorithm constructs a max-min fair allocation for a network, N . In plain English, the algorithm iterates over a set of receivers, each step increasing those receivers' rates uniformly as much as possible without overutilizing any links in the network. A receiver is removed from this set once some link on its data-path reaches full capacity, or, if the receiver is part of a single-rate session, the data-path of some receiver in the session contains a link that has reached full capacity. We define

$$\phi_{i,j}(T) = \begin{cases} 1 & |R_{i,j} \cap T| > 0 \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

1. $T_0 = \{r_{i,k}\}; \forall r_{i,k}, a_{i,k}^0 = 0; \forall i, j, u_{i,j}^0 = 0, u_j^0 = 0; b = 0$
2. While $|T_b| > 0$

3. $t_{b+1} = \sup\{t : \forall j, u_j^b + \sum_i \phi_{i,j}(T_b)t \leq c_j \wedge \forall r_{i,k} \in T_b \Rightarrow a_{i,k}^b + t \leq \alpha_i\}$
4. $\forall r_{i,k} \in T_b, a_{i,k}^{b+1} = a_{i,k}^b + t_{b+1}$. For all other $r_{i,k}, a_{i,k}^{b+1} = a_{i,k}^b$.
5. $u_{i,j}^{b+1} = \sum_{r_{i,k} \in R_j} a_{i,k}^{b+1}, u_j^{b+1} = \sum_i u_{i,j}^{b+1}$.
6. $T' = T_b - \{r_{i,k} \in T_b : a_{i,k}^{b+1} = \alpha_i \vee (\exists j, r_{i,k} \in R_{i,j} \wedge u_j^{b+1} = c_j)\}$
7. $T_{b+1} = T' - \{r_{i,k} \in T' : \sigma(S_i) = S \wedge \exists r_{i,k'} \notin T'\}$
8. $b++$
9. end while
10. $\forall r_{i,k}, a_{i,k} = a_{i,k}^b, \forall i, j, u_{i,j} = u_{i,j}^b, u_j = u_j^b$

Step 3 the largest value that all receivers' rates in T_b can be incremented while maintaining feasibility of the allocation. Steps 4 and 5 apply this increase to the "current" receiver rates and link rates respectively. Step 6 removes any receivers from T_{b+1} whose rates cannot be incremented any further, or else they would be larger than the maximum session rate, or would cause overutilization of some link. Step 7 removes any receivers in single-rate sessions from T_{b+1} , given that some other receiver in that session has been removed (so that all receiver rates in this session remain identical).

Lemma 5 *The above algorithm's resulting allocation is max-min fair.*

Proof: The choice of t_{b+1} in step 3 causes some link to be fully utilized, or causes some receiver to attain its session's maximum rate, α_i so that $T_b \supset T_{b+1}, T_b \neq T_{b+1}$, hence the algorithm must terminate. This choice also ensures that $u_j^{b+1} \leq c_j$ for each b , hence the final allocation must be feasible. Since $\forall b, t_b > 0, a_{i,k} > a_{i',k'} \iff \exists \beta, r_{i,k} \in T_\beta$ and $r_{i',k'} \notin T_\beta$. It follows that $a_{i,k} = a_{i',k'} \iff \exists \beta, r_{i,k}, r_{i',k'} \in T_\beta$ and $r_{i,k}, r_{i',k'} \notin T_{\beta+1}$.

Let A represent the allocation produced by the algorithm, and let \bar{A} be another feasible allocation (using $\bar{a}_{i',k'}$ to represent each receiver's rate in \bar{A}). Consider any receiver $r_{i,k}$ where $\bar{a}_{i,k} > a_{i,k}$. To show that A is max-min fair, we must show that there exists some other receiver $r_{i',k'}$ where $\bar{a}_{i',k'} < a_{i',k'} \leq a_{i,k}$.

Let β represent the iteration of the algorithm (i.e., the value of b) where $r_{i,k} \in T_\beta$, but $r_{i,k} \notin T_{\beta+1}$. If $r_{i,k}$ is excluded from $T_{\beta+1}$ as a result of step 7, then S_i is in a single-rate session, and some other receiver $r_{i,k'}$ in the same session must have been removed in step 6. Because S_i is single-rate, we must have that $a_{i,k} = a_{i,k'}$ and $\bar{a}_{i,k} = \bar{a}_{i,k'}$ and thus $\bar{a}_{i,k'} = \bar{a}_{i,k} > a_{i,k} = a_{i,k'}$. Hence, there is some receiver that was excluded from $T_{\beta+1}$ as a result of step 6 whose rate allocations in A and \bar{A} are identical to those of $r_{i,k}$. Hence, we can assume, WLOG, that $r_{i,k}$ was excluded as a result of step 6. Because of this, and because $a_{i,k} < \bar{a}_{i,k} \leq \alpha_i$, it must be that $r_{i,k}$'s data-path contains a link l_j that is fully utilized at the end of iteration β . Thus, any other receiver $r_{i',k'} \in R_j$ is excluded from T_{b+1} , so that $a_{i',k'} \leq a_{i,k}$. Since l_j is fully utilized and $\bar{a}_{i,k} > a_{i,k}$ there is some receiver $r_{i',k'}$ where $\bar{a}_{i',k'} < a_{i',k'}$ to prevent l_j from being utilized beyond its capacity. Since this $r_{i',k'} \in R_j$, we have $\bar{a}_{i',k'} < a_{i',k'} \leq a_{i,k}$. ■

C Multi-rate Max-Min Fairness Proofs

We now present proofs for several of the Lemmas and Theorems in the paper. To ensure that there is no circular usage of proofs (e.g., Lemma A's proof uses Lemma B, whose proof uses Lemma A), each proof uses only the lemmas, theorems and corollaries that are *proven* previously. For instance, we do not establish that the max-min fair allocation is unique until Corollary 5, so all proofs prior must assume that more than one max-min fair allocation may exist.

Lemma 6 *A multi-rate max-min fair allocation is fully-utilized-receiver-fair.*

We prove that a max-min fair allocation is fully-utilized-receiver-fair by showing for each receiver $r_{i,k}$, its rate $a_{i,k}$ in a max-min fair allocation A is fully-utilized-receiver-fair. We begin by arbitrarily choosing the receiver $r_{i,k}$. If $a_{i,k} = \alpha_i$, then the proof holds trivially. Since $a_{i,k} \leq \alpha_i$ in any feasible allocation, we need only consider when $a_{i,k} < \alpha_i$. Let $\gamma = \min_{\{j: r_{i,k} \in R_j\}} c_j - u_j$ (i.e., an amount of bandwidth that is free on all links on $r_{i,k}$'s data-path). Since A is feasible, $\gamma \geq 0$. We construct another allocation, \bar{A} , as follows. $\bar{a}_{i',k'} = a_{i',k'}$ for any receiver $r_{i',k'} \neq r_{i,k}$, and $\bar{a}_{i,k} = a_{i,k} + \gamma$. Then \bar{A} is feasible, but if $\gamma > 0$, then $\bar{a}_{i,k} > a_{i,k}$ and there is no receiver $r_{i',k'}$ where

$\bar{a}_{i',k'} < a_{i',k'} \leq a_{i,k}$, contradicting max-min fairness of A . Hence, we must have $\gamma = 0$, which means there must be some fully utilized link l_j for which $r_{i,k} \in R_j$.

It remains to show that there is some link l_j that is fully utilized where $r_{i,k} \in R_j$ and $a_{i',k'} \leq a_{i,k}$ for all other $r_{i',k'} \in R_j$. For all links l_j , define $\beta_j = \max_{r_{i',k'} \in R_j} \{0, a_{i',k'} - a_{i,k}\}$, and let $\beta = \min_{r_{i,k} \in R_j \wedge u_j = c_j} \beta_j$ (β is well-defined since we have shown some link on $r_{i,k}$'s data-path must be fully utilized). Also define $\epsilon = \min_{r_{i,k} \in R_j \wedge u_j < c_j} \{\beta, c_j - u_j\}$. Since each $\beta_j \geq 0$ and $c_j - u_j \geq 0$ in a feasible allocation, we have that $\epsilon \geq 0$. Note from its definition that $\epsilon = 0$ only if $\beta = 0$.

We will now increase and decrease some receiver rates by ϵ to produce a feasible allocation that contradicts the max-min fairness of A (unless $\epsilon = 0$). Let \bar{A} be this new allocation where $\bar{a}_{i,k} = a_{i,k} + \epsilon$, and $\bar{a}_{i',k'} = a_{i',k'}$ whenever $a_{i',k'} \leq a_{i,k} + \epsilon$, and $\bar{a}_{i',k'} = a_{i',k'} - \epsilon > a_{i,k}$ otherwise. If a receiver's allocation is smaller in \bar{A} than in A , then its allocation in \bar{A} is bounded below by $a_{i,k}$, so all receiver allocations in \bar{A} are non-negative. We write \bar{u}_j for the link rate of l_j under allocation \bar{A} . Since only receiver $r_{i,k}$'s allocation is larger in \bar{A} than in A , it follows that for any link l_j , if $r_{i,k} \notin R_j$ then $0 \leq \bar{u}_j \leq u_j \leq c_j$. For $r_{i,k} \in R_j$, if $\epsilon = 0$, clearly $\bar{u}_j = u_j$. If instead $\epsilon > 0$, for $u_j < c_j$, the construction of ϵ gives that $\epsilon \leq c_j - u_j$, and since only $r_{i,k}$'s allocation is increased in \bar{A} , we have $\bar{u}_j \leq u_j + \epsilon \leq c_j$. If $u_j = c_j$, then, by definition of β_j , there is at least one receiver $r_{i',k'} \in R_j$ where $a_{i',k'} - a_{i,k} = \beta_j \geq \beta \geq \epsilon > 0$. Having $r_{i',k'} \in R_j$ ensures that $\bar{u}_j \leq u_j \leq c_j$, thus \bar{A} is feasible (no links are utilized beyond capacity).

If $\epsilon > 0$, then $\bar{a}_{i,k} > a_{i,k}$, and for each $r_{i',k'}$ where $\bar{a}_{i',k'} < a_{i',k'}$, we have that $a_{i',k'} > a_{i,k}$. Hence, \bar{A} 's feasibility contradicts the max-min fairness of A . Thus, $\epsilon = 0$, and it follows that $\beta = 0$. Thus, some $\beta_{j'} = 0$ where $r_{i,k} \in R_{j'}$. The definition of β_j gives us that $u_{j'} = c_{j'}$, and for all other $r_{i',k'} \in R_{j'}$, $a_{i',k'} - a_{i,k} \leq 0$. Hence, $a_{i,k}$ is fully-utilized-receiver-fair.

Since $r_{i,k}$ was chosen arbitrarily, each receiver's rate is fully-utilized-receiver-fair, making the allocation A fully-utilized-receiver-fair. ■

Corollary 2 *A multi-rate max-min fair allocation is shared-path-receiver-fair.*

Proof: Consider any pair of receivers $r_{i,k}$ and $r_{i',k'}$ (i may or may not equal i') that have the same data path. By Lemma 6, there exists a link l_j where $\forall r_{i'',k''} \in R_j, a_{i,k} \geq a_{i'',k''}$, hence $a_{i,k} \geq a_{i',k'}$. Lemma 6 also gives us that there exists a link $l_{j'}$ (perhaps even the same link) where $\forall r_{i'',k''} \in R_{j'}, a_{i',k'} \geq a_{i'',k''}$, hence $a_{i',k'} \geq a_{i,k}$. It follows that $a_{i',k'} = a_{i,k}$. ■

Lemma 7 *In a multi-rate max-min fair allocation, for each receiver $r_{i,k}$, either $a_{i,k} = \alpha_i$, or else there is at least one fully utilized link, l_j , where $r_{i,k} \in R_{i,j}$ and for all sessions $S_{i'}, u_{i',j} \leq a_{i,k} \leq u_{i,j}$.*

Proof: If $a_{i,k} = \alpha_i$, then the proof holds trivially. Assume $a_{i,k} < \alpha_i$. Then by Lemma 6, there is a fully utilized link l_j on $r_{i,k}$'s path from the sender where $a_{i',k'} \leq a_{i,k}$ for all receivers $r_{i',k'} \in R_j$. Hence $u_{i',j} = \max\{a_{i',k'} | r_{i',k'} \in R_{i',j}\} \leq a_{i,k}$ for each session, $S_{i'}$. In Section 2, we defined $u_{i,j}$ to be $u_{i,j} = \max\{a_{i,k'} | r_{i,k'} \in R_{i,j}\}$, making $u_{i,j} \geq a_{i,k}$. ■

Corollary 3 *A multi-rate max-min fair allocation is per-receiver-link-fair.*

Proof: Follows directly from Lemma 7.

Corollary 4 *A multi-rate max-min fair allocation is per-session-link-fair.*

Proof: This follows easily from Corollary 3.

Proof of Theorem 1: The proof is the immediate result of Lemma 6, Corollary 2, Lemma 7, and Corollaries 3 and 4. ■

Lemma 8 *Let S_i be a single-rate session in a network N and $a_i < \alpha_i$. In a max-min fair allocation, there exists a fully utilized link l_j where $|R_{i,j}| > 0$ and $a_i = u_{i,j} \geq u_{i',j}$ for all other sessions $S_{i'}$.*

Proof: Let A be the max-min fair allocation. Let $\gamma = \min_{\{j: |R_{i,j}| > 0\}} c_j - u_j$. Since A is feasible, $\gamma \geq 0$. Let \bar{A} be an allocation where $\bar{a}_{i',k'} = a_{i',k'}$ for any receiver $r_{i',k'} \notin S_i$, and $\bar{a}_{i,k} = a_{i,k} + \gamma$ for all $r_{i,k} \in S_i$. Note that all receivers in S_i are reduced by an identical amount, so all receivers in S_i continue to receive at the same rate. Since the rate is increased only in session S_i , a link l_j 's link rate increases by at most γ , and the increase only occurs when

$|R_{i,j}| > 0$, so by the definition of γ , the link rate remains beneath c_j . Then \bar{A} is feasible, and for the same reasons as in Lemma 6, contradicts the max-min fairness of allocation A unless $\gamma = 0$. Hence, there must be some fully utilized link l_j for which $r_{i,k} \in R_j$ for some receiver $r_{i,k} \in S_i$.

Let $\beta_j = \max_{r_{i',k'} \in R_j} \{0, a_{i',k'} - a_i\}$, and let $\beta = \min_{|R_{i,j}| > 0 \wedge u_j = c_j} \beta_j$. Also define $\epsilon = \min_{|R_{i,j}| > 0 \wedge u_j < c_j} \{\beta, c_j - u_j\}$. Then let \bar{A} be an allocation where $\bar{a}_i = a_i + \epsilon$, and $\bar{a}_{i',k'} = a_{i',k'}$ whenever $a_{i',k'} \leq a_{i,k} + \epsilon$, and $\bar{a}_{i',k'} = a_{i',k'} - \epsilon$ otherwise. Using similar reasoning as in the second half of the proof of Lemma 6, some $\beta_{j'}$ must be 0, so that on link l_j , $a_i = u_{i,j} \geq u_{i',j}$ for all other sessions $S_{i'}$. ■

Proof of Lemma 1: For each receiver $r_{i,k}$, let $a_{i,k}$ be the receiver's rate in allocation A , and $\bar{a}_{i,k}$ be the receiver's rate in allocation B . We also write $A = (\gamma_1, \dots, \gamma_s)$ as the ordered vector of receiver rates in A and $B = (\beta_1, \dots, \beta_s)$ as the ordered vector of receiver rates in B .

If $B = A$ then $B \dashv A$. Otherwise, some receiver's rate must differ in allocations A and B . If all receiver rates are lower in B than in A , then clearly $B \dashv A$. Now consider the case where there is at least one receiver where $\bar{a}_{i,k} > a_{i,k}$. Let $M = \{r_{i,k} : \bar{a}_{i,k} > a_{i,k}\}$, and choose a (possibly unique) receiver r_{i_1,k_1} from M with the minimal value of $a_{i,k}$, ($a_{i_1,k_1} = \min_{r_{i,k} \in M} a_{i,k}$). Since A is a max-min solution, and since $a_{i_1,k_1} < \bar{a}_{i_1,k_1} \leq \alpha_{i_1}$, by Lemma 6 if $\sigma(S_{i_1}) = \mathcal{M}$, or by Lemma 8 if $\sigma(S_{i_1}) = \mathcal{S}$, there is a link l_j that is fully utilized where $r_{i_1,k_1} \in R_j$. To prevent link l_j from being overutilized in allocation B , there must be some other receiver $r_{i_2,k_2} \in R_j$ where $a_{i_2,k_2} > \bar{a}_{i_2,k_2}$. Furthermore, Lemmas 6, 8 give us that $a_{i_1,k_1} \geq a_{i,k}$ for all other receivers $r_{i,k} \in R_j$, hence $a_{i_1,k_1} \geq a_{i_2,k_2} > \bar{a}_{i_2,k_2}$.

Return now to our ordered vectors of A and B . Consider the smallest d' where $\gamma_{d'} = a_{i_2,k_2}$ (such a d' must exist, since r_{i_2,k_2} is allocated a_{i_2,k_2} in A). If $d' = 1$, then since β_1 equals the lowest rate assigned to any receiver in allocation B , making $\beta_1 \leq \bar{a}_{i_2,k_2} < a_{i_2,k_2} = \gamma_1$, and we have $B \dashv A$. If $d' > 1$, then (by the choice of d') $\gamma_{d'-1} < a_{i_2,k_2}$, which means there are $d' - 1$ receivers, each of whose rate $a_{i,k}$ satisfies $a_{i,k} < a_{i_2,k_2} < a_{i_1,k_1}$. From our choice of r_{i_1,k_1} within M , any receiver $r_{i,k}$ that satisfies $a_{i,k} < a_{i_1,k_1}$ must also satisfy $\bar{a}_{i,k} \leq a_{i,k}$. Thus we have that each of these $d' - 1$ receivers has $\bar{a}_{i,k} \leq a_{i,k}$, and since each of these receiver's rates in allocation A satisfies $a_{i,k} < a_{i_2,k_2} = \gamma_{d'}$, we have identified a set X of $d' - 1$ receivers whose rates in allocation B are less than $\gamma_{d'}$. Since each receiver $r_{i,k} \in X$ satisfies $a_{i,k} < a_{i_2,k_2}$, receiver $r_{i_2,k_2} \notin X$ (since its allocated rate in A is a_{i_2,k_2} , and its allocation in B is $\bar{a}_{i_2,k_2} < a_{i_2,k_2}$). Thus, including r_{i_2,k_2} , there are at least d' receivers whose allocation in B is less than a_{i_2,k_2} so $\beta_{d'} < a_{i_2,k_2} = \gamma_{d'}$.

For any d where $\gamma_d < a_{i_1,k_1}$, there are at least d receivers whose allocations in A each satisfy $a_{i,k} < a_{i_1,k_1}$, and hence (from the choice of r_{i_1,k_1}) each of these receiver's rates satisfies $\bar{a}_{i,k} \leq a_{i,k}$. Since there are also at least d receivers that satisfy $a_{i,k} \leq \gamma_d$, we have that $\bar{a}_{i,k} \leq a_{i,k} \leq \gamma_d$ for at least d receivers, which means that there are at least d receivers whose rates in allocation B are less than γ_d , hence $\beta_d \leq \gamma_d$. Since (by choice of d') $\gamma_{d'-1} < a_{i_2,k_2} \leq a_{i_1,k_1}$, it follows that $\gamma_d < a_{i_1,k_1}$ for $d < d'$, thus $\beta_d \leq \gamma_d$ for $d < d'$. We therefore have $\beta_d \leq \gamma_d$ for $d < d'$ and $\beta_{d'} < \gamma_{d'}$. ■

Note that the proof makes no assumptions about the type of sessions in the network (unicast, multi-rate single-rate multicast). Thus, the proof holds for a network with any combinations of types of sessions.

Proof of Lemma 2: Let $X = (x_1, \dots, x_k), Y = (y_1, \dots, y_k)$.

If: Since $X \neq Y, X \dashv Y$, there must be at least one index i where $x_i < y_i$ (otherwise all $x_i \geq y_i$, and since $X \neq Y$, for some $i, x_i > y_i, X \not\vdash Y$). Let $i_{x'} = \min_i \{i : x_i < y_i\}$, and let $x' = x_{i_{x'}}$, so that $|\{x_i : x_i \leq x'\}| \geq i_{x'}$. Since $y_{i_{x'}} > x_{i_{x'}} = x'$, we have that $i_{x'} > |\{y_i : y_i \leq x'\}|$. Thus, $|\{x_i : x_i \leq x'\}| > |\{y_i : y_i \leq x'\}|$.

Any $j < i_{x'}$ implies that $x_j \leq x_{i_{x'}} = x'$, so that $x_j = y_j$. Otherwise, if $x_j > y_j$, since $X \dashv Y$, there must be some $i' < j$ where $x_{i'} < y_{i'}$, contradicting our choice of $i_{x'}$. Similarly, $x_j < y_j$ would contradict our choice of $i_{x'}$. Since $x_j = y_j$ for all $j < i_{x'}$, clearly for any $z < x', |\{x_i \in X : x_i \leq z\}| = |\{y_i \in Y : y_i \leq z\}|$, and it follows trivially that for any $z < x', |\{x_i \in X : x_i \leq z\}| \geq |\{y_i \in Y : y_i \leq z\}|$.

Only If: Define $i_{x'} = |\{x_i \in X : x_i \leq x'\}|$. Since by assumption, $|\{y_i \in Y : y_i \leq x'\}| < |\{x_i \in X : x_i \leq x'\}| = i_{x'}$, we have $x_{i_{x'}} \leq x' < y_{i_{x'}}$. This assumption also gives us that if $y_j = x'$, then $x_j \leq x'$, hence $y_j = x'$ implies that $y_j \geq x_j$. It is also the case that when $y_j < x'$, then $y_j \geq x_j$. To prove this, assume instead that $y_j < x_j$. This means that, $|\{x_i \in X : x_i \leq y_j\}| < j$. Furthermore, it is always the case that $j \leq |\{y_i \in Y : y_i \leq y_j\}|$. Thus, $|\{x_i \in X : x_i \leq y_j\}| < |\{y_i \in Y : y_i \leq y_j\}|$ for $y_j < x'$, which contradicts the assumption that for all $z < x', |\{x_i \in X : x_i \leq z\}| \geq |\{y_i \in Y : y_i \leq z\}|$. To summarize, $y_{i_{x'}} > x' \geq x_{i_{x'}}$. For $j < i_{x'}$, if $y_j \geq x'$,

because $x_j \leq x_{i,w} \leq x'$, we have $y_j \geq x_j$. Furthermore, if $y_j < x'$, then $y_j \geq x_j$. It follows that $X \dashv Y, X \neq Y$. ■

Corollary 5 (Max-min fair uniqueness) *There is a unique max-min fair allocation for any network.*

Proof: Let A and B be max-min fair allocations for a network, N . Then by Lemma 1, $A \dashv B$ and $B \dashv A$, hence $A = B$. For each receiver $r_{i,k}$, let $a_{i,k}$ be its rate under allocation A , and $b_{i,k}$ be its rate under allocation B . Note that equality of the ordered vectors does not imply that $a_{i,k} = b_{i,k}$ for each receiver $r_{i,k}$. To prove this, define $\gamma = \min_{r_{i,k}} \{a_{i,k} : a_{i,k} \neq b_{i,k}\}$. Note that because $A = B$, our choice of γ gives us that for any receiver $r_{i,k}$ where $a_{i,k} = \gamma$, then $a_{i,k} \leq b_{i,k}$ (this follows from Lemma 2 where $x' = \gamma$).

Fix $r_{i,k}$ to be any (possibly unique) receiver $r_{i,k}$ where $b_{i,k} > a_{i,k} = \gamma$ (one such receiver exists given the construction of γ and the fact that $\gamma = a_{i,k} \neq b_{i,k} \rightarrow a_{i,k} \leq b_{i,k}$). Since B is max-min fair, it is feasible, and from our choice of $r_{i,k}$, there is no receiver $r_{i',k'}$ where $b_{i',k'} < a_{i',k'} \leq a_{i,k}$. This contradicts the max-min fairness of A . Thus, for each receiver $r_{i,k}$ within the network, we must have $a_{i,k} = b_{i,k}$. ■

Proof of Theorem 2: Since we have that the max-min fair allocation is unique, and that the algorithm in Appendix B computes the max-min fair allocation, we can simply examine the allocation computed by the algorithm.

(a) A receiver $r_{i,k} \in S_i$ where $\sigma(S_i) = \mathcal{M}$ is removed in step 6. Hence, $a_{i,k} = \alpha_i$, or there is some fully utilized link l_j that led to $r_{i,k}$'s exclusion from $T_{\beta+1}$ for some β . Following the argument in Lemma 5, $a_{i,k}$ is the largest possible rate for any receivers whose data-path crosses l_j .

(b) Apply the argument in Lemma 7, replacing Lemma 6 with (a)

(c) Per-session-link-fairness holds for each session S_i where $\sigma(S_i) = \mathcal{M}$ as a consequence of (b). For $\sigma(S_i) = \mathcal{S}$, at least one receiver $r_{i,k}$ must be excluded from $T_{\beta+1}$ for some β by step 6. If $a_{i,k} = \alpha_i$, then since S_i is single-rate, $a_{i,k'} = a_{i,k} = \alpha_i$ for all $r_{i,k'} \in S_i$. Otherwise, there is some link l_j that is fully utilized and, again following the argument in Lemma 5, $a_{i,k}$ is the largest possible rate for any receivers whose data-path crosses l_j . Since for all sessions $S_{i'}$, we have that $u_{i',j} = \max\{a_{i',j} : r_{i',j} \in R_{i',j}\}$, it follows that $u_{i,j} \geq u_{i',j}$ for all other $S_{i'}$.

(d) follows from (a) using an argument similar to that in the proof of Corollary 2. (e) follows from directly from (a). ■

Lemma 9 *If a single session S_i switches from being single-rate to multi-rate, its max-min fair rates per receiver do not decrease.*

Proof: Let N be the network where $\sigma(S_i) = \mathcal{M}$, and \bar{N} the network where $\sigma(S_i) = \mathcal{S}$, let A be the max-min fair allocation in N , and \bar{A} the max-min fair allocation in \bar{N} . For an arbitrary receiver $r_{i',k}$, let $\bar{a}_{i',k}$ be its max-min fair rate in \bar{N} , and $a_{i',k}$ be its max-min fair rate in N . For $r_{i,k} \in S_i$, Theorem 2(a) gives us, there exists fully utilized link l_j where $a_{i,k} = \alpha_i$ or $a_{i,k} \geq a_{i',k'}$ for all $r_{i',k'} \in R_j$. If $\bar{a}_{i,k} > a_{i,k}$, then to prevent l_j from being utilized beyond capacity, there must be some other receiver $r_{i',k'} \in R_j$ where $\bar{a}_{i',k'} < a_{i',k'}$. However, since $a_{i,k} < \bar{a}_{i,k} \leq \alpha_i$, $a_{i,k} \geq a_{i',k'}$, so we have $\bar{a}_{i,k} \geq a_{i,k} \geq \bar{a}_{i',k'} \geq a_{i',k'}$. Since \bar{A} is a feasible allocation in N , this results in a contradiction that A is the max-min fair allocation. Thus, it must be the case for each $r_{i,k} \in S_i$ that $\bar{a}_{i,k} \leq a_{i,k}$. ■

D Additional Properties of max-min fair allocations

We now expand on our discussion of the properties of max-min fair allocations. In Section 2, we showed how removal of a receiver from a session can cause other receivers' multi-rate max-min fair rates to vary in either direction (increase or decrease), regardless of whether or not the receiver is in the same session as the removed receiver. We now present results that prove that max-min fair receiver rates do not drop below the max-min fair rate of the receiver that is removed, and hence, a single-rate session's rate never decreases when a receiver is removed from a session.

Lemma 10 *Let $N = \{G, \{S_1, \dots, S_m\}, \tau, \sigma\}$ and $\bar{N} = \{G, \{\bar{S}_1, S_2, \dots, S_m\}, \tau, \sigma\}$ be networks which differ only in that session S_1 has one additional receiver, r_{1,k_1} , than session \bar{S}_1 . For each receiver $r_{i,k}$, let $a_{i,k}$ be its max-min fair rate in network N , and (except for receiver r_{1,k_1} , let $\bar{a}_{i,k}$ be its max-min fair rate in network \bar{N} . If $a_{i,k} < a_{1,k_1}$, then $\bar{a}_{i,k} = a_{i,k}$, and if $a_{i,k} = a_{1,k_1}$, then $\bar{a}_{i,k} \geq a_{i,k}$.*

Proof: Consider any receiver $r_{i,k}$ where $a_{i,k} < a_{1,k_1}$, and apply the algorithm in Appendix B which constructs the max-min fair allocation to both networks N and \bar{N} . When applying the algorithm to N , since $a_{i,k} < a_{1,k_1}$, there

exists β such that $r_{i,k} \notin T_\beta$ but $r_{1,k_1} \in T_\beta$. Thus, no link whose full-utilization would prevent r_{1,k_1} 's inclusion in T_b can be fully utilized whenever $b \leq \beta$. Furthermore, $\alpha_1 > \sum_{i=1}^{\beta} t_i$. Thus, the allocation to r_{1,k_1} does not effect the value of $a_{i',k'}^t$ for any other receiver $r_{i',k'}$ where $t \leq \beta$. Thus, for $t \leq \beta$, $a_{i',k'}^t$ is identical over N and \bar{N} . Since $r_{i,k}$ is excluded from T_β , we have that $\bar{a}_{i,k} = a_{i,k} = a_{i,k}^\beta$.

We now show that if $a_{i,k} = a_{1,k_1}$, then $\bar{a}_{i,k} \geq a_{i,k}$. Let $A = (\alpha_1, \dots, \alpha_s)$ be the ordered vector of the allocation of rates to receivers in \bar{N} , where each receiver $r_{i,k}$'s allocation equals $a_{i,k}$ (i.e., each receiver except r_{1,k_1} receives at the rate which is its max-min fair rate in N). Clearly, this is a feasible allocation for \bar{N} , so by Lemma 1, $A \dashv \bar{A}$, where $\bar{A} = (\bar{\alpha}_1, \dots, \bar{\alpha}_s)$ is the ordered vector of max-min fair receiver rates in \bar{N} . If $A = \bar{A}$, then $n = |\{a_{i,k} \leq a_{1,k_1}\}| = |\{\bar{a}_{i,k} \leq a_{1,k_1}\}|$, otherwise $\alpha_n \neq \bar{\alpha}_n$, contradicting $A = \bar{A}$. If $A \neq \bar{A}$, then by Lemma 2, $\exists x'$ such that $\forall z \leq x', |\{\alpha_i \in A : \alpha_i \leq z\}| \geq |\{\bar{\alpha}_i \in \bar{A} : \bar{\alpha}_i \leq z\}|$ and $|\{\alpha_i \in A : \alpha_i \leq x'\}| > |\{\bar{\alpha}_i \in \bar{A} : \bar{\alpha}_i \leq x'\}|$. Since all receivers where $a_{i',k'} < a_{1,k_1}$ implies $\bar{a}_{i',k'} = a_{i',k'}$, we have that $x' \geq a_{1,k_1}$. This means that if there is any receiver $r_{i',k'}$ where $a_{i',k'} \geq a_{1,k_1} > \bar{a}_{i',k'}$, then $|\{\bar{\alpha}_i \in \bar{A} : \bar{\alpha}_i \leq \bar{a}_{i',k'}\}| = |\{\alpha_i \in A : \alpha_i \leq \bar{a}_{i',k'}\}| + 1$, which is a contradiction. Thus, $\bar{a}_{i,k} \geq a_{i,k} = a_{1,k_1}$. ■

Clearly, Lemma 10 can be extended to any session S_i , i.e., i need not equal 1. The Lemma states that a receiver $r_{i,k}$'s max-min fair rate will not change due to the removal of a receiver from the network when $a_{i,k}$ is less than the removed receiver's max-min fair rate (prior to its removal). Furthermore, the max-min fair rate for a receiver $r_{i,k}$ can only increase by removing a receiver when the two receiver rates are equal, prior to the receiver removal. This gives an interesting result about single-rate session rates due to the removal of a receiver.

Corollary 6 *Let N be a network. If a session S_i is single-rate, then if a receiver leaves the session, the session's max-min fair rate can only increase.*

Proof: All receivers $r_{i,k} \in S_i$ have identical rates. From Lemma 10, by removing a receiver, r_{i,k_i} from the session, since we have $a_{i,k} \leq a_{i,k_i}$ prior to the removal (in fact, they are equal), hence $\bar{a}_{i,k} \geq a_{i,k}$. ■

Corollary 6 gives us the following important result: a single-rate session's max-min fair rate does not decrease (but may increase) when a receiver is removed from within the session.

E One and Two Layer Expected Bandwidth with Random Joins

We compute the expected bandwidth for session S_i on a link l_j . For simplicity, we write $R = |R_{i,j}|$, and denote the set of receivers from session S_i whose data-path utilizes this link (i.e., $R_{i,j}$) as $\{r_1, \dots, r_R\}$, and let a_t be the number of packets that receiver r_t must receive per quantum.

Let ρ packets be transmitted in a time quantum, and let X_i be a random variable that equals 1 if any receiver is joined when packet i is transmitted, and 0 otherwise ($1 \leq i \leq \rho$). Let $Y_{i,t}$ be a random variable that equals 1 if receiver r_t joins to receive packet i , and 0 otherwise. Since we assume a receiver chooses the packets it is to receive from a uniform distribution, we have $\Pr(Y_{i,t} = 1) = a_t/\rho$.

$$\begin{aligned} E[X_i] &= 1 - \prod_{t=1}^R \Pr(Y_{i,t} = 0) = 1 - \prod_{t=1}^R (1 - a_t/\rho) \\ E[U_{i,j}] &= E\left[\sum_{i=1}^{\rho} X_i\right] = \sum_{i=1}^{\rho} E[X_i] = \rho \left(1 - \prod_{t=1}^R (1 - a_t/\rho)\right) \end{aligned}$$

Let us now consider the inclusion of a second layer. For simplicity, we present the formula when all receivers require the same number of packets per quantum, $\forall t, a_t = a$. If f is the fraction of packets transmitted on the bottom layer, if $a \geq f$, then any packet transmitted on the bottom layer has $E[X_i] = 1$, and for a packet on the top layer, $\Pr(Y_{i,t} = 1) = \frac{a-f\rho}{(1-f)\rho}$. When $a < f$, any packet transmitted on the top layer has $E[X_i] = 0$, and for a packet on the bottom layer, $\Pr(Y_{i,t} = 1) = \frac{a}{f\rho}$. This yields:

$$E[U_{i,j}] = \left\{ \begin{array}{ll} \rho(f + (1-f)(1 - \left(1 - \frac{a-f\rho}{(1-f)\rho}\right)^R)) & a \geq f \\ \rho f(1 - (1 - \frac{a}{f\rho})^R) & a < f \end{array} \right\} = \left\{ \begin{array}{ll} \rho(f + (1-f)(1 - \left(\frac{1-a/\rho}{1-f}\right)^R)) & a \geq f \\ \rho f(1 - (1 - a/f\rho)^R) & a < f \end{array} \right\}$$

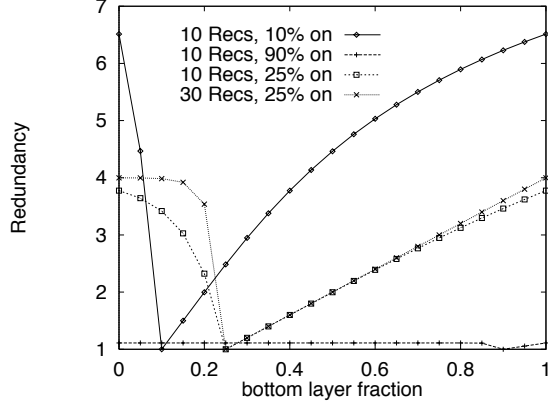


Figure 9: Redundancy of a single layer with random joins

Let us consider the benefits of using multiple layers. Figure 9 examines how having two layers reduces redundancy compared to having a single layer. The aggregate rate of transmission by the session is split between the two layers, where the x -axis indicates the fraction, f , of packets that are transmitted on the bottom layer per quantum. Any receiver whose fair rate is $a_{i,k} \geq \rho f$ remains joined to the lower layer and randomly selects the remaining packets per quantum from the upper layer. Any receiver whose rate is $a_{i,k} < \rho f$ randomly chooses its $a_{i,k}/\rho$ packets per quantum off the lower layer. Each curve represents a fixed number of receivers, each of which receives an identical percentage of packets per quantum (from Figure 5, we know that identical percentages yield the highest redundancy). When $f = 0$ and $f = 1$, all packets are transmitted on a single layer, hence the redundancy is identical to when there is only one layer available. For $0 < f < 1$, redundancy is lower than in the single layer case, and equals exactly one when f and the percentage of packets needed per quantum by the receiver with the largest receiving rate are equal.

Using layers in this manner reduces the randomness in the packet selection process by a receiver. This increases the correlation of packets chosen, hence having an additional layer decreases redundancy. However, since the maximum percentage of packets received per quantum can vary from link to link within a session, it is difficult to choose an optimal value for f .

F Details of Congestion Control Approaches

The congestion control protocols used in Section 4 are based mainly on the ideas in [18], but vary in several respects in order to increase the “history-less” nature of the protocol, simplifying its Markov model. Most of these differences are discussed in Section 4. The one additional difference is that the sender does not transmit packets at a fixed rate on each of the layers. Instead, it transmits on layer i at an *expected rate* of 1 for layer 1, and at 2^{i-2} for each layer $i > 1$. Hence, the expected rate for a receiver joined up to layer i is simply $1 + \sum_{j=2}^i 2^{j-2} = 2^{i-1}$. Hence, the expected rates for the protocols used in Section 4 and the protocol presented in [18] are identical.

Let us now describe how the sender achieves an expected rate of 1 on layer 1 and of 2^{i-2} on layer i where $i > 1$, when the sender transmits over N layers. The sender transmits packets over all layers at rate 2^{N-1} . It places a packet on layer i with probability 2^{i-1-N} for $i > 1$, and with probability 2^{1-N} on layer 1.

F.1 Implementing the Coordinated Protocol

We now discuss implementation details of the Coordinated protocol. For a protocol that uses N layers, a packet sent on layer 1 requires a $\log_2 N$ bit field to implement the coordination. Receivers join to an additional layer only when receiving a packet transmitted on layer 1. Upon receipt of such a packet, receivers examine the value in this field. If the value of the field, j , is larger than the layer i up to which the receiver is joined ($j > i$), then the receiver joins layer $i + 1$. The probability that the field is set to the value $j < N - 1$ is $1/2^j$. The probability that the field is set to the value $N - 1$ is $1/2^{N-2}$. The conditional probability that a receiver, joined up to layer i , receives a packet sent on layer 1 given that it receives a packet on some layer is $1/2^{i-1}$. The conditional probability that the packet causes the receiver to join an additional layer, given that the packet arrives on layer 1 is $\sum_{j=i}^{N-2} 1/2^j + 1/2^{N-2} = 1/2^{i-1}$. Hence,

the conditional probability that the receiver joined up to layer i joins an additional layer upon receiving a packet is $1/2^{i-1} * 1/2^{i-1} = 1/2^{2(i-1)}$ as desired.

F.2 Markov Model of 2-receiver session

We now describe the Markov models used to compute the expected data rates of the Uncoordinated and Coordinated protocols for a session containing two receivers, configured in a modified-star topology (Figure 7(a)). The models give the expected data rate on the shared link, as well as on each of the two fanout links. We first describe the models used in the Uncoordinated and Coordinated protocols, since these models are quite similar. We then extend the model to a Semi-Coordinated protocol; motivation for a Semi-Coordinated protocol is discussed there as well.

Consider a session that transmits over N layers. We model such a session for Coordinated and Uncoordinated protocols using N^2 states, where each state is labeled (i, j) , $1 \leq i, j \leq N$. The state labeled (i, j) represents a session configuration where receiver 1 is joined up to layer i , and receiver 2 is joined up to layer j . Define $\rho(i)$ to be the expected rate at which data is transmitted on layers 1 through layer i . When in state (i, j) , receiver 1 receives at expected rate $\rho(i)$, receiver 2 at expected rate $\rho(j)$, and the shared link transmits data at expected rate $\rho(\max\{i, j\})$.

A transition is taken from state (i, j) whenever a packet is transmitted from the sender on any layer from 1 to $\max\{i, j\}$. We call such transmission a *transmission event*. After a transmission event, the new current state of the Markov process is (i', j') where $i' = i - 1$ if receiver 1 loses the packet, $i' = i$ if receiver 1 does not lose the packet, or gets the packet and decides not to add an additional layer, or $i' = i + 1$ if receiver 1 receives the packet and decides to add an additional layer. Receiver 2 behaves in identical manner to determine whether $j' = j - 1, j$, or $j + 1$.

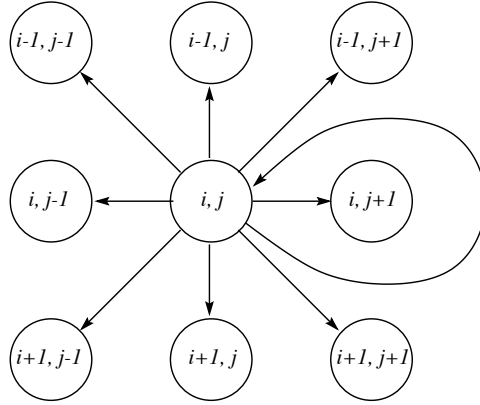


Figure 10: A typical state and the transitions from that state in the 2 receiver Markov model.

Figure 10 illustrates the transitions that can occur for a typical state (i, j) . States that are not typical are those where $i = 1, N$ or $j = 1, N$. In these cases, states (i', j') where $i' = 0, N + 1$ or $j' = 0, N + 1$ and connecting transitions are simply omitted.

The models for the Coordinated protocol and the Uncoordinated protocol assign different weights to the directed transitions between pairs of states. Assuming the Markov process is ergodic, we can compute the steady state probabilities, $\Pi_{i,j}$ of residing within state (i, j) . $\Pi_{i,j}$ is the steady-state probability, receiver 1 is joined up to layer i and receiver 2 is joined up to layer j when a transmission event occurs.

Let $T_{(i,j)}$ be a random variable indicating the time to a subsequent transmission event upon arriving in state (i, j) , and define $\mathcal{O}_{(i,j)}(o)$ to be a random variable that equals 1 whenever some observer o is joined up to the layer on which the subsequent transmission event is sent, and 0 otherwise. Let R_o be a random variable that equals the rate at which observer o obtains packets. Then

$$E[R_o] = \sum_{1 \leq i, j \leq N} \Pi_{i,j} E[\mathcal{O}_{(i,j)}(o)] / \sum_{1 \leq i, j \leq N} \Pi_{i,j} E[T_{(i,j)}] \quad (2)$$

For each transmission by the sender, we define λ_k to be the event that a packet transmitted by the sender is sent on some layer l , $1 \leq l \leq k$. Upon transitioning into state (i, j) , the shared link carries the subsequent transmission event

with probability 1, receiver 1 will observe the subsequent transmission event with probability $\Pr(\lambda_i|\lambda_{\max\{i,j\}})$, and receiver 2 observes the this transmission event with probability $\Pr(\lambda_j|\lambda_{\max\{i,j\}})$.

Given that the sender sends packets at a rate of $\rho(N)$, the expected time until a subsequent transmission in state (i, j)

$$E[T_{(i,j)}] = \frac{1}{\rho(N)\Pr(\lambda_i)}. \quad (3)$$

We are interested in the expected rates of three ‘‘observers’’: the shared link, receiver 1, and receiver 2. Define $\mathcal{O}_{(i,j)}(s)$, $\mathcal{O}_{(i,j)}(1)$, and $\mathcal{O}_{(i,j)}(2)$ to be random variables that, upon entering state (i, j) , equal 1 or 0 depending on whether the shared link, receiver 1, and receiver 2 (respectively) are joined up to the layer on which the subsequent transmission event is sent. Our definition of a transmission event gives us the following:

$$E[\mathcal{O}_{(i,j)}(s)] = 1 \quad (4)$$

$$E[\mathcal{O}_{(i,j)}(1)] = \begin{cases} 1 & i \geq j \\ \Pr(\lambda_i|\lambda_j) & i < j \end{cases} \quad (5)$$

$$E[\mathcal{O}_{(i,j)}(2)] = \begin{cases} 1 & i \leq j \\ \Pr(\lambda_j|\lambda_i) & i > j \end{cases} \quad (6)$$

We now compute transition weights for the set of states in the Markov model. Let us consider a state (i, j) and assume that $i \geq j$. The transition from this state when a transmission event occurs depends on the outcome of three events. First, we must consider whether or not receiver 2 is joined up to the layer on which the packet is transmitted (we make the reverse consideration for receiver 1 when $j > i$). Next, for each receiver that is joined up to the layer on which the packet is transmitted, we must determine whether the packet is dropped, either on the shared link or on individual link. Finally, for each receiver that receives the packet, we must consider whether the packet causes the receiver to join an additional layer. For this purpose, we define several events that are used in the calculations of transition weights for all protocols. Recall that λ_j is the event that receiver 2 (joined up to layer j) is joined to the layer on which the transmission event is transmitted. Let G_1 be the event that the path to receiver 1 is not congested. When a transmission event is sent on a layer joined by receiver 1, if G_1 holds, then the packet is received, otherwise it is lost. We define G_2 in a similar manner for receiver 2. Let I_k be the event that a receiver joined up to layer k joins an additional layer if it receives a transmission event. Note that as defined, the event G_1 is independent of λ_i, λ_j , and I_k . The same can be said for G_2 . Also, I_k is independent of both λ_i and λ_j .

Table 1: State transitions and their weights for a state (i, j) , $1 < j \leq i < N$.

| transition | value |
|---------------------------------|---|
| $(i, j) \rightarrow (i-1, j-1)$ | $\Pr(\lambda_j)\Pr(\neg G_1 \wedge \neg G_2)$ |
| $(i, j) \rightarrow (i-1, j)$ | $\Pr(\lambda_j)\Pr(\neg G_1 \wedge G_2)\Pr(\neg I_j) + \Pr(\neg \lambda_j)\Pr(\neg G_1)$ |
| $(i, j) \rightarrow (i-1, j+1)$ | $\Pr(\lambda_j)\Pr(\neg G_1 \wedge G_2)\Pr(I_j)$ |
| $(i, j) \rightarrow (i, j-1)$ | $\Pr(\lambda_j)\Pr(G_1 \wedge \neg G_2)\Pr(\neg I_i)$ |
| $(i, j) \rightarrow (i, j)$ | $\Pr(\lambda_j)\Pr(G_1 \wedge G_2)\Pr(\neg I_i \wedge \neg I_j) + \Pr(\neg \lambda_j)\Pr(G_1)\Pr(\neg I_i)$ |
| $(i, j) \rightarrow (i, j+1)$ | $\Pr(\lambda_j)\Pr(G_1 \wedge G_2)\Pr(\neg I_i \wedge I_j)$ |
| $(i, j) \rightarrow (i+1, j-1)$ | $\Pr(\lambda_j)\Pr(G_1 \wedge \neg G_2)\Pr(I_i)$ |
| $(i, j) \rightarrow (i+1, j)$ | $\Pr(\lambda_j)\Pr(G_1 \wedge G_2)\Pr(I_i \wedge \neg I_j) + \Pr(\neg \lambda_j)\Pr(G_1)\Pr(I_i)$ |
| $(i, j) \rightarrow (i+1, j+1)$ | $\Pr(\lambda_j)\Pr(G_1 \wedge G_2)\Pr(I_i \wedge I_j)$ |

Table 1 gives the values for the weights for transitions from state (i, j) where $1 < j \leq i < N$, in terms of the probabilities of the events defined above. Weights for transitions of states (i, j) where $1 < i < j < N$ can be computed in a similar manner (i.e., in that case λ_j holds for each transmission event, whereas λ_i may or may not hold). Note that because $j \leq i$, we have $\Pr(\lambda_j) = 1$. Hence, we need not include terms that contain $\neg \lambda_i$. It also follows that $\Pr(\lambda_j|\lambda_i) = \Pr(\lambda_j)$, and $\Pr(\neg \lambda_j|\lambda_i) = \Pr(\neg \lambda_j)$.

Table 2: Probabilities for join correlations for state $(i, j), 1 \leq j \leq i \leq N$.

| Event | uncorrelated value | correlated value |
|---------------------------------|-------------------------------|-----------------------|
| $\Pr(I_i \wedge I_j)$ | $\Pr(I_i) \Pr(I_j)$ | $\Pr(I_i)$ |
| $\Pr(\neg I_i \wedge I_j)$ | $\Pr(\neg I_i) \Pr(I_j)$ | $\Pr(I_j) - \Pr(I_i)$ |
| $\Pr(I_i \wedge \neg I_j)$ | $\Pr(I_i) \Pr(\neg I_j)$ | 0 |
| $\Pr(\neg I_i \wedge \neg I_j)$ | $\Pr(\neg I_i) \Pr(\neg I_j)$ | $\Pr(\neg I_j)$ |

Transition weights leading from states where $i = 1, N$ or $j = 1, N$ must be modified slightly: receivers never drop layer 0, nor can they add a layer beyond N . The easiest way to compute the weights for transitions from such states can be done using the entries in Table 1: We demonstrate for the transition $(1, N) \rightarrow (1, N)$, whose weight is the sum of table entries $(1, N) \rightarrow (1, N)$, $(1, N) \rightarrow (1, N + 1)$, $(1, N) \rightarrow (0, N)$, and $(1, N) \rightarrow (0, N + 1)$.

The Coordinated protocol and the Uncoordinated protocol differ only in the dependence relation between I_i (the join event for receiver 1) and I_j (the join event for receiver 2). For the Uncoordinated protocol, these events are independent. For the Coordinated protocol, recall that if $i \geq j$, then event I_i holds only if I_j holds as well. We find that for both protocols, for each state $(i, j), 1 \leq j \leq i \leq N$, the probabilities of various combinations of events I_i and I_j holding and not holding can be written as functions of $\Pr(I_i)$ and $\Pr(I_j)$. These relations are given in Table 2.

We conclude by listing the values of the other probability parameters in Table 1 for an arbitrary state (i, j) .

- $\Pr(\lambda_j) = \max\{1, \rho(i)/\rho(j)\}$.
- $\Pr(\neg G_1 \wedge \neg G_2) = p_s + (1 - p_s)p_1p_2$
- $\Pr(\neg G_1 \wedge G_2) = (1 - p_s)p_1(1 - p_2)$
- $\Pr(G_1 \wedge \neg G_2) = (1 - p_s)(1 - p_1)p_2$
- $\Pr(G_1 \wedge G_2) = (1 - p_s)(1 - p_1)(1 - p_2)$
- $\Pr(G_1) = (1 - p_s)(1 - p_1)$
- $\Pr(G_2) = (1 - p_s)(1 - p_2)$

$\Pr(I_i)$ and $\Pr(I_j)$ are defined explicitly for both protocols to be $1/2^{2(i-1)}$ and $1/2^{2(j-1)}$ respectively. $\rho(k)$ is defined to be 2^{k-1} .

F.3 Extending the model to the Semi-Coordinated protocol

We now extend our Markov model to include the Semi-Coordinated protocol. The Semi-Coordinated protocol behaves similarly to the Uncoordinated protocol, except when receivers are joined to an identical layer after losing the same packet. When this occurs, receiver joins are synchronized until one receiver drops a packet that the other does not. The Semi-Coordinated protocol captures a feature of the Deterministic protocol used in the simulations in Section 4: when two receivers join or leave the same layer at the same time, then all subsequent joins will occur at the same time, until a loss event causes only one of the receivers to leave a layer.

We introduce N additional states into the Markov model that we label $(0, 1)$ through $(0, N)$, where state $(0, i)$ represents a session in which the two receivers both are joined up to layer i and the prior loss observed by each is the same.

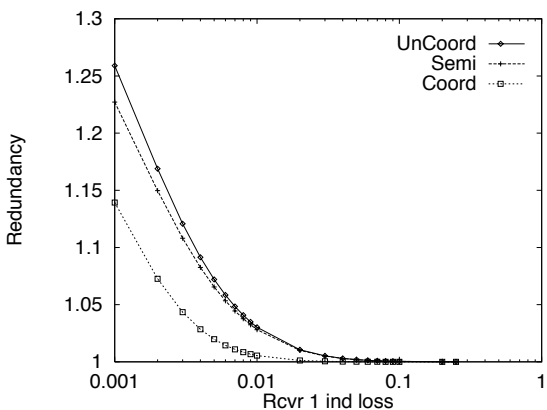
The transitions for any state $(i, j) \rightarrow (i', j')$ for $1 \leq i, j \leq N$ are identical to those in the Markov model for the Uncoordinated protocol, with the exception of those listed in Table 3.

F.4 Markov Model results

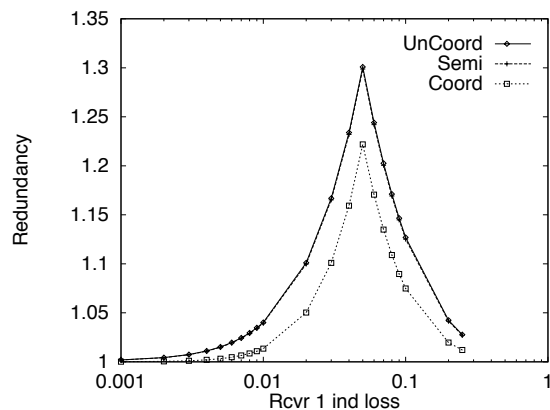
The protocols operating with two receivers over seven layers are evaluated via the Markov models. Figures 11 and 12 compare the redundancies of the two protocols under a variety of loss conditions. In each graph, we vary p_1 along the x -axis. The y -axis gives the redundancy, curves represent the results for the Coordinated, Uncoordinated, and Semi-Coordinated protocols. This Semi-Coordinated protocol attempts to capture differences between our Uncoordinated protocol, and such a protocol where join events occur after a deterministic number of packets. In particular, the Semi-Coordinated protocol captures the fact that when two receivers join or leave the same layer at the same time, then all

Table 3: Transition differences between Semi-Coordinated and Coordinated

| Transition | Semi-Coordinated Value |
|--|---|
| $(1, 1) \rightarrow (1, 1)$ | $\Pr(G_1 \wedge G_2) \Pr(\neg I_1) \Pr(\neg I_2) + (\Pr(\neg G_1 \wedge G_2) + \Pr(G_1 \wedge \neg G_2)) \Pr(\neg I_1)$ |
| $(i, i) \rightarrow (i-1, i-1), 1 < i \leq N$ | 0 |
| $(i, i) \rightarrow (0, i-1), 1 < i \leq N$ | $\Pr(\neg G_1 \wedge \neg G_2)$ |
| $(0, 1) \rightarrow (0, 1)$ | $\Pr(G_1 \wedge G_2) \Pr(\neg I_1) + \Pr(\neg G_1 \wedge \neg G_2)$ |
| $(0, i) \rightarrow (0, i), 1 < i \leq N$ | $\Pr(G_1 \wedge G_2)$ |
| $(0, i) \rightarrow (0, i+1), 1 \leq i < N$ | $\Pr(G_1 \wedge G_2) \Pr(I_i)$ |
| $(0, N) \rightarrow (0, N)$ | $\Pr(G_1 \wedge G_2)$ |
| $(0, 1) \rightarrow (1, 1)$ | $(\Pr(\neg G_1 \wedge G_2) + \Pr(G_1 \wedge \neg G_2)) \Pr(\neg I_1)$ |
| $(0, i) \rightarrow (i', j'), 1 \leq i \leq N$ $1 \leq i', j' \leq N, i' \neq j', i' - i \leq 1, j' - j \leq 1$ | same as Coordinated model's $(i, i) \rightarrow (i', j')$ |
| All other transitions are identical to the coordinated model | |



(a) Low independent R2 loss (.001)



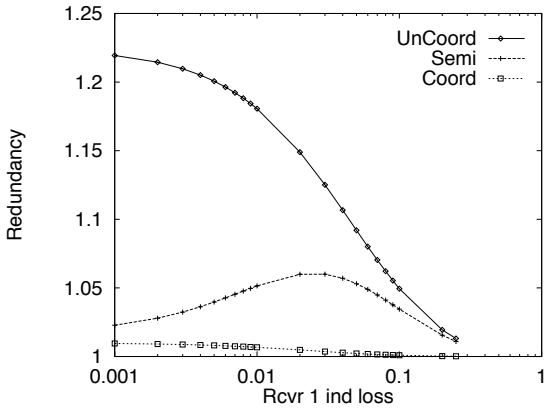
(b) High independent R2 loss (.05).

Figure 11: Low shared loss (.001), 7 layers

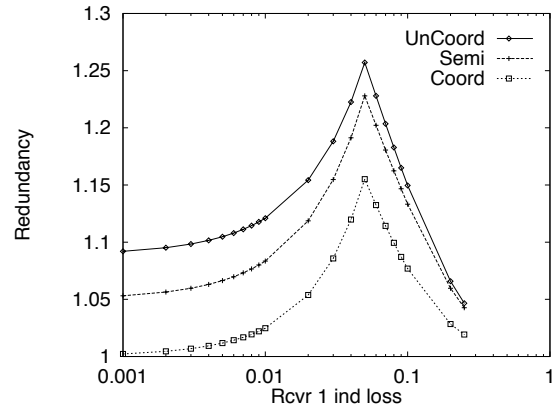
subsequent joins will occur at the same time, until a loss event causes only one of the receivers to leave a layer. It is possible, but pointless to implement the Semi-Coordinated protocol.

We consider four combinations of a low loss rate (.001) and a high loss rate (.05) for p_s and p_2 . From the figures, we conclude that the highest redundancy occurs when the independent loss rates are equal. This is depicted by the peaks in Figures 11(b) and 12(b). The peak is due to the change in the receiver whose rate is used to calculate the optimum expected utilization of the link (the higher rate always belongs to the receiver with higher loss). High shared loss causes the greatest variation in Coordinated and Uncoordinated protocols: when shared loss is high, receivers are more likely to drop layers at the same time, so the difference in redundancy due to the differences in sender coordination is more pronounced. The Coordinated protocol's redundancy tends to be smaller when at least one receiver's independent loss is very low. The Semi-Coordinated protocol's redundancy is most often very close to the redundancy of the Uncoordinated protocol, unless shared loss is very high and independent losses are very low, or there is a significant difference in receivers' independent losses. Hence, for most network loss scenarios, whether or not the join period occurs after a random point in time or fixed point in time does not influence redundancy, as long as the expected time to join is the same.

The most important point we draw from these graphs is that redundancy is mostly affected by the variation in receiver loss rates, and is highest when these loss rates are identical.



(a) Low independent R2 loss (.001)



(b) High independent R2 loss (.05)

Figure 12: High shared loss (.05), 7 layers