

Smooth Multirate Multicast Congestion Control

Gu-In Kwon
guin@cs.bu.edu

John W. Byers
byers@cs.bu.edu

Computer Science Department
Boston University
Boston, MA 02215

Abstract—A significant impediment to deployment of multicast services is the daunting technical complexity of developing, testing and validating congestion control protocols fit for wide-area deployment. Protocols such as pgmcc and TFMCC have recently made considerable progress on the single rate case, i.e. where one dynamic reception rate is maintained for all receivers in the session. However, these protocols have limited applicability, since scaling to session sizes beyond tens of participants necessitates the use of multiple rate protocols. Unfortunately, while existing multiple rate protocols exhibit better scalability, they are both less mature than single rate protocols and suffer from high complexity.

We propose a new approach to multiple rate congestion control that leverages proven single rate congestion control methods by orchestrating an ensemble of independently controlled single rate sessions. We describe SMCC, a new multiple rate equation-based congestion control algorithm for layered multicast sessions that employs TFMCC as the primary underlying control mechanism for each layer. SMCC combines the benefits of TFMCC (smooth rate control, equation-based TCP friendliness) with the scalability and flexibility of multiple rates to provide a sound multiple rate multicast congestion control policy.

I. INTRODUCTION

Despite considerable effort and numerous technical advances, a suitable multiple rate multicast congestion control mechanism fit for wide area deployment is still yet to emerge. The primary reason appears to be the daunting complexity associated with delivering different TCP-friendly rates to different participants within the session. In all existing schemes for multiple rate congestion control, versions of layered multicast (originally proposed in [1]) are employed, whereby different multicast groups within the multicast session transmit at different rates, and participants use IGMP messages to join and leave groups to adjust their rate. But the significant challenges associated with this method are that the actions of one receiver can adversely impact other receivers; moreover, joins can place sudden load on the network, leading to unfriendliness to protocols such as TCP. Existing methods to mitigate these problems ultimately lead to very complex multiple rate congestion control designs that are difficult to evaluate.

Further evidence of the technical hurdles associated with multiple rate schemes is given by promising recent advances in *single rate* multicast congestion control, notably pgmcc [2] and TFMCC [3]. With single rate congestion control schemes, the sender transmits at a rate requested by the slowest

receiver in the group. While these protocols are not designed to scale to large or heterogeneous audiences, there is building consensus that these protocols are sufficiently mature and well-tested for Internet deployment. In this paper, we seek to leverage these advances. In particular, we explore a new direction in multiple rate multicast congestion control, namely building a multiple rate scheme from an ensemble of single rate sessions, *each of which has their own independent control*. The major advantage of this method is that it leverages proven single rate congestion control mechanisms to provide an effective multiple rate scheme with relatively little additional complexity. This is in contrast to all existing multiple rate congestion control schemes, which provide only an integrated control mechanism *across* layers, and do not attempt to take advantage of control mechanisms *within* layers. As a result, these integrated controls are often extremely complex, and are difficult to test and validate.

A. Our Work in Context

There has been a significant amount of previous work on TCP-friendly multiple rate multicast congestion control, including [4], [5], [6], [7], [1], [8]. All of these approaches employ layered multicast, i.e. they employ a set of multicast groups that transmit at different rates to accommodate a heterogeneous, and potentially large population of receivers. Previous work has categorized these schemes as either using static or dynamic layers. In static schemes, such as [1], [7], the sending rate of any given layer remains fixed over time, and all adjustments to the reception rate are therefore exclusively receiver-driven. This approach has some drawbacks, most notably that the receiver may have insufficient information to accurately conduct join attempts, as well as necessitating abrupt rate changes. Many other schemes use dynamic layers, or layers whose transmission rate changes over time according to a predetermined pattern. Dynamic layers have been used in a variety of clever ways, including implicit coordination of receivers behind a bottleneck [8], reduction of IGMP leave messages [4], simulation of additive increase [6], and to achieve equation-based congestion control [9]. However, implementations of these dynamic layering schemes typically have a great deal of embedded complexity to realize these benefits in practice.

One feature shared by most existing multiple rate methods is that the layer rates are *non-adaptive*, i.e. the schedule of packet

transmissions on each group (whether fixed-rate or dynamic) is known to the sender and to the receivers in advance. A limitation of non-adaptive schemes is their inflexibility; there are typically only a small constant number of feasible control actions that may be taken by a receiver at a given time step. For example, in many non-adaptive schemes, the magnitude by which a receiver may instantaneously increase or decrease its rate is fixed a priori, and the times at which a rate increase can be performed are widely separated. Our work differs in this regard, since each of the TFMCC sessions comprising the individual layers adaptively and continuously adjust their rates to the limiting receivers in the session, as we will describe.

Two methods for adaptive, layered multiple rate multicast were proposed in SAMM [10] and HALM [11]. However, the methods proposed in SAMM predated current notions of TCP-friendliness and were not evaluated in that context, moreover, extra router support to monitor the available bandwidth is required to achieve the best performance. The work in HALM is most similar to our own, in that they advocate periodic, adaptive reallocation of layer rates in a multirate multicast session and build upon single rate congestion control methods. In their case though, the emphasis is on periodic optimization of the layer rates at coarse time scales (tens of seconds) that is not suitable for fine-grained congestion control on the Internet.

B. Contributions and Organization

We describe SMCC, a new multiple rate equation-based congestion control algorithm for cumulative layered multicast sessions that employs TFMCC as the primary underlying control mechanism for each layer. Since each layer is controlled independently by TFMCC, the properties of TFMCC hold for participants in any given layer. As such, the layer rates are both dynamic and adaptive. SMCC combines the benefits of TFMCC (smooth rate control, equation-based TCP friendliness) with the scalability and flexibility of multiple rates to provide a sound multiple rate multicast congestion control policy. In SMCC, each receiver cumulatively subscribes to appropriate layers based on its estimated rate using the TCP throughput equation [12] also employed in TFMCC. In addition to the TFMCC functionality, SMCC provides a new additive increase join attempt to avoid abrupt rate increases when the receiver attempts to subscribe to an additional layer. Ultimately, the smooth rate change of SMCC is ideally suited to streaming multimedia applications; but equation-based methods are general-purpose, thus SMCC works naturally for other multicast applications, such as reliable downloads [8], [13]. Finally, it is worth emphasizing that SMCC requires no additional router support beyond basic multicast functionality, and does not place any new demands on any existing multicast protocols.

The remainder of this paper is organized as follows. In Section II, we review the underlying TFMCC congestion control mechanism. In Section III, we specify SMCC and how to orchestrate an ensemble of TFMCC sessions to build a multirate congestion control scheme. In Section IV, we propose a new additive increase join attempt which is performed by each

receiver before joining the next layer. In Section V, we give the results from ns simulations to demonstrate the fairness of SMCC with competing TCP flows.

II. TFMCC OVERVIEW

TFMCC [3] is a single rate multicast congestion control protocol designed to provide smooth rate change over time. TFMCC extends the basic equation-based control mechanisms of TFRC [14] into the multicast domain. The fundamental idea is to have each receiver evaluate a control equation (Eqn. 1) derived from the model of TCP's long-term throughput [12], then use this to directly control the sender's transmission rate.

$$T_{TCP} = \frac{s}{RTT \left(\sqrt{\frac{2p}{3}} + (12\sqrt{\frac{3p}{8}})p(1 + 32p^2) \right)} \quad (1)$$

where T_{TCP} is a function of the steady-state loss event rate p , the TCP round-trip time RTT , and the packet size s .

A cursory overview of TFMCC functionality is as follows:

- Each receiver measures the packet loss rate.
- The receiver measures or estimates the round-trip time to the sender.
- The receiver uses the control equation (Eqn. 1) to derive an acceptable transmission rate from the measured loss rate and round-trip time.
- The receiver sends the calculated transmission rate to the sender.
- A feedback suppression scheme (additional details below) is used to prevent feedback implosion while ensuring that feedback from the slowest receiver always reaches the sender.
- The sender adjusts the sending rate from the feedback information.

In TFMCC, the receiver that the sender believes currently has the lowest expected throughput of the group is selected as the *current limiting receiver* (CLR). The CLR sends continuous, immediate feedback to the sender without any suppression, so the sender can use the CLR's feedback to adjust the transmission rate. In addition, any receiver whose expected throughput is lower than the sender's current rate sends a feedback message, and to avoid feedback implosion, biased feedback timers in favor of receivers with lower rates are used.

A. Measuring the Loss Event Rate

One crucial detail of TFMCC which we will return to later in the paper is the method it uses to measure packet loss. In TFMCC, a receiver aggregates the packet losses into *loss events*, defined as one or more packets lost during a round-trip time. The number of packets between consecutive loss event is called a *loss interval*. The average loss interval size can be computed as the weighted average of the m most recent loss intervals l_k, \dots, l_{k-m+1} :

$$l_{avg}(k) = \frac{\sum_{i=0}^{m-1} w_i l_{k-i}}{\sum_{i=0}^{m-1} w_i}$$

The weights w_i are chosen so that very recent loss intervals receive the same high weights, while the weights gradually decrease to 0 for older loss intervals. The loss event rate p used as an input for the TCP model is then taken to be the inverse of l_{avg} . The interval since the most recent loss event is incomplete, since it does not end with a loss event, but it is conservatively included in the calculation of the loss event rate if doing so reduces p :

$$p = \frac{1}{\max(l_{avg}(k), l_{avg}(k-1))}$$

B. Round-trip Time Measurements

Each receiver starts with an initial RTT estimate that is used until a real measurement is made. A receiver measures the RTT by sending timestamped feedback to the sender, which then echoes the timestamp and receiver ID in the header of a data packet. An exponentially weighted moving average (EWMA) is used to prevent a single large RTT measurement from greatly impacting the sending rate.

$$t_{RTT} = \beta \cdot t_{RTT}^{inst} + (1 - \beta) \cdot t_{RTT}$$

The recommended value of β for the CLR is 0.05 while all other receivers are recommended to use $\beta = 0.5$ due to their less infrequent RTT measurements. For further details of TFRC and TFMCC, we refer the reader to [14] and [3].

III. SMOOTH MULTIRATE MULTICAST CONGESTION CONTROL

Like many other multiple rate congestion control schemes, SMCC employs cumulative layered multicast. But unlike other schemes, SMCC employs adaptive layers, i.e. each individual layer uses TFMCC congestion control and each receiver subscribes to appropriate layers based on its calculated equation-based rate. The high-level features of our approach are as follows:

- Each receiver subscribes to a set of cumulative layers. We refer to a receiver as being an *active* participant in the uppermost layer to which it subscribes, and a *passive* participant in all other layers.
- Each layer i of SMCC transmits at a rate within a designated interval and the rate floats within that interval according to TFMCC congestion control regulated by active participants in that layer.
- The current limiting receiver (CLR) for each layer is selected from among the active participants of that layer to adjust the sending rate.
- Each receiver calculates its expected throughput.
- If the expected throughput calculated from the equation is above the maximum sending rate of its current subscription level, the receiver performs a join attempt using additive increase methods.
- If a receiver's computed throughput is below the minimum receiving rate of the layer i , it drops its highest layer i . (Note that this bounds the extent to which a CLR can drag down a single TFMCC session).

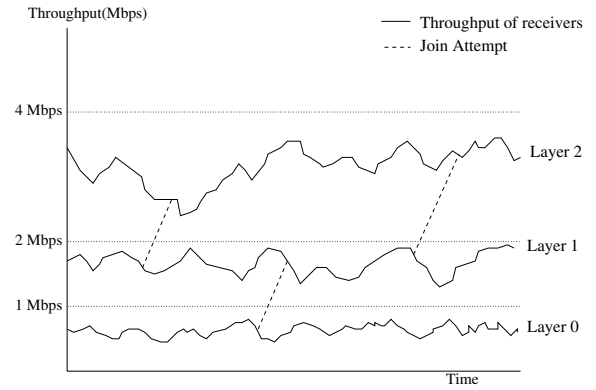


Fig. 1. SMCC Overview

In the following section, we describe how to set up the layers and define how the CLR is selected for each layer. Then, in Section III-B, we describe how each receiver sends feedback and how changes of the CLR on each layer are realized. In Section IV, we describe our methods for additive increase join attempts.

Term	Description
B_i	the maximum cumulative sending rate up through layer i
C_i	CLR on layer i
L_i	layer i
R_i	actual sending rate on layer i
S_j	subscription level of receiver j
T_i	aggregate target rate requested by C_i

TABLE I
SMCC TERMS

A. Setting up Layers and CLR

We employ a cumulative layering scheme so that each receiver subscribes and unsubscribes to layers in sequential order. For simplicity, in the following discussion and in the remainder of the paper, we will assume that the maximum cumulative sending rates through layer i , which we denote by B_i , follow the natural 1, 2, 4, 8, ... progression. Figure 1 briefly shows the configuration of layers and how the sending rate of each layer is set, and Table I describes the terms we use in the following sections to describe SMCC. Our approach is amenable to other multiplicative layer rate increases, as advocated in [4], or to finer-grained rates of increase. We define the maximum cumulative sending rates of the layers formally as follows: We let B_0 be the maximum sending rate of the base layer, and we set $B_i = 2^i * B_0$ for $i \geq 1$. From this setting of the rates, we can associate each desired reception rate with a set of subscription layers: a receiver j desiring rate r_j should subscribe to all layers i such that $B_i \leq 2r_j$. In addition, a receiver which has a computed throughput in the range $[0, B_0]$ always subscribes to the base layer L_0 . In this sense, we can map each receiver to the layer on which they are active. We say that layer L_i is *responsible* for receivers with rates in the range $[B_{i-1}, B_i]$. Equivalently,

we define the *subscription level* S_j of receiver j to be the layer responsible for that receiver. The subscription level of receiver with expected throughput x is:

$$S_j = \left\lceil \log_2 \frac{x}{B_0} \right\rceil.$$

For example, a receiver with expected throughput 6 Mbps where $B_0 = 1$ Mbps has a subscription level of 3 (i.e. it subscribes to L_0, L_1, L_2 , and L_3) and L_3 is responsible for this receiver. At any instant in time, we let C_i denote the current limiting receiver of a given layer i , i.e. the active receiver j that has the lowest expected throughput in the range $[B_{i-1}, B_i]$.

B. Dynamically Adjusting Layer Rates

Now we consider the setting of the layer rates, starting with the base layer, L_0 . The sender adjusts the sending rate of the base layer based on the feedback sent by C_0 , the CLR for the base layer, and we denote the actual sending rate on the base layer that results from this process by R_0 . The other active receivers on the base layer do not send feedback unless their calculated rate is less than R_0 . In all other respects, the sender and the active receivers in L_0 follow the TFMCC scheme.

Receivers with expected throughput in the range of $[B_0, 2B_0]$ subscribe to L_1 as well as L_0 . Let T_1 denote the total aggregate rate requested by the limiting receiver subscribing to L_1 . Then, the actual sending rate R_1 on layer 1 is set to the *difference* between T_1 and R_0 .

In general, the same principle is used to set the rate R_i on layer i :

$$R_i = T_i - \sum_{j=0}^{i-1} R_j, \quad (2)$$

where T_i is the aggregate target rate requested by C_i , the current limiting receiver on L_i . From this setting, it is easy to show the following bounds on R_i :

$$\forall i : 0 \leq R_i \leq B_i - B_{i-2}.$$

As in the TFMCC approach, the active participants in L_i do not send feedback unless their calculated rate is less than T_i , thus avoiding feedback implosion. The CLRs are permitted to send immediate feedback without any form of suppression, so the sender can use the CLRs' feedback to adjust the transmission rate (upward or downward) for each layer.

The CLR for a layer can change in one of two ways: either a new receiver becomes the CLR or the existing CLR leaves the group. Each of these cases is relatively easy to handle. If a receiver whose subscription level is i sends feedback that indicates a rate that is lower than the current rate of C_i , but still larger than B_{i-1} , the sender will set C_i to that receiver and immediately reduces its rate for L_i to the requested rate in the feedback message according to Equation (2). If a receiver on L_i has a calculated rate which is less than B_{i-1} , it unsubscribes from layer L_i . The receiver needs to issue one IGMP leave message to drop the layer. While dropping the highest layer does not guarantee a particular

amount of multiplicative decrease, on average, the reception rate is decreased by half.

If the departing receiver is the CLR on L_i , a new CLR for layer i must be elected. To accomplish this, a departing CLR first sends a control message to the sender notifying it of the departure. Upon receipt of this signal, the sender multicasts a control message to the group asking active participants to select a new CLR. As in TFMCC, each receiver which is an active participant on layer L_i will set a random timer before sending feedback to the sender. To avoid feedback implosion, biased feedback timers in favor of receivers with lower rates are used.

If there are no active participants on layer i (which can happen when other participants are active on other layers j such that $j > i$), no CLR is assigned to layer i . The actual sending rate of layer i is then set to $(B_i - \sum_{j=0}^{i-1} R_j)$ and the rates on higher layers are adjusted according to Equation (2). If any receiver which is active in layer $j > i$ subsequently drops layers $i + 1$ through j and becomes active in layer i , this receiver will become the CLR in layer i , as will a receiver who joins layer i from below. The sending rate of layer i is then adjusted by this active receiver's feedback rate.

C. Subscribing to an Additional Layer

Even though the receivers in the same group have similar calculated throughput, they may not share the same congested links. So, measured packet loss events across active receivers in a layer will vary. Often, some receivers may compute a calculated throughput value which is in the range of the next layer, and those receivers will attempt to join the next layer. As motivated in the introduction and in related work [1], naive join attempts using a single IGMP join request are problematic, as they introduce a sudden rate increase along a network path. Such a spurious join attempt may cause significant packet loss prior to the time at which the attempt is rescinded [4]. In severe cases, this substantial increase on the bottleneck link may drive TCP flows into timeout. For this reason, join attempts which mimic fine-grained additive increase [5], [6] are preferable. Here, instead of joining the next layer, the receiver increases the receiving rate slowly, i.e. by one more packet per RTT, during the join attempt.

Another compelling reason for proceeding to the next layer slowly is due to inaccuracies in estimating the target throughput when it differs substantially from the current reception rate using TFMCC methods. As described earlier in section II-A, the loss rate is computed from the loss interval, which is defined as the number of received packets since the last loss event. Hence, the loss interval clearly depends on the sending rate. But since the sending rate is controlled by the CLR's feedback, the loss rate currently measured by a non-CLR is not the same as if the sending rate adjusted to its feedback. In section V, we show simulation results demonstrating that the loss rate measured by non-CLR is not a sufficiently accurate estimate to conclusively determine whether or not to join the next layer. In practice, depending on the specific

scenario considered, the calculated throughput can either be an overestimate or an underestimate.

Our methods for performing additive increase joins are the glue that holds an ensemble of TFMCC sessions together, and constitute the key additional feature needed to provide a sound multiple rate congestion control scheme. As such, we describe them fully in Section IV.

D. Avoiding Oscillation

Recall that both loss rate and RTT are input parameters to calculate the transmission rate, thus even slight increases in RTT cause the calculated rate to be slightly reduced. If we allow the receiver to join layer L_{i+1} when the receiver's calculated transmission rate is very close to B_i , even a slight increase of RTT is enough to force the receiver to drop the highest layer so that it goes back to its previous layer. This may cause significant oscillation when the RTT fluctuates.

We employ a conservative method to avoid this oscillation. The receiver will join the next layer L_{i+1} if the calculated rate is in the range of $[\alpha \cdot B_i, B_{i+1}]$ for a *damping factor* α such that $1 < \alpha < 2$. A fixed value, such as $\alpha = 1.2$, can be used for all receivers. In practice, methods for computing a receiver-specific damping factor which incorporates the variance in that receiver's RTT can produce better performance and less oscillation. We describe those methods more fully in the full version of the paper [15].

IV. ADDITIVE INCREASE JOIN ATTEMPTS

We now describe a new additive increase scheme to conduct join attempts between successive layers in our multicast session. Although other work has proposed the use of additive increase in multiple rate multicast congestion control, such as FGLM [5] and STAIR [6], those methods are designed as an integral part of complex, non-cumulative multicast layering schemes, and have technical limitations which make them unsuitable for this application. In contrast, the layers we propose for additive increase are *only* used when a receiver wishes to attempt to join the next successive layer. Our scheme has the following properties.

- True additive increase with respect to end-to-end bandwidth consumption.
- Employs no IGMP messages (which can be slow to take effect).
- Uses only a small number of additional IGMP join messages.

Once a receiver performing a join attempt from layer L_i attains a total reception rate equal to T_{i+1} , the target rate sent by C_{i+1} , it joins layer L_{i+1} and drops the special additive increase layers. If, however, there is a packet loss during the join attempt, the receiver ceases the join attempt. We incorporate the information gained from both successful or failed join attempts into loss interval and loss rate calculations. The sender sends the next layer rate information in the packet header.

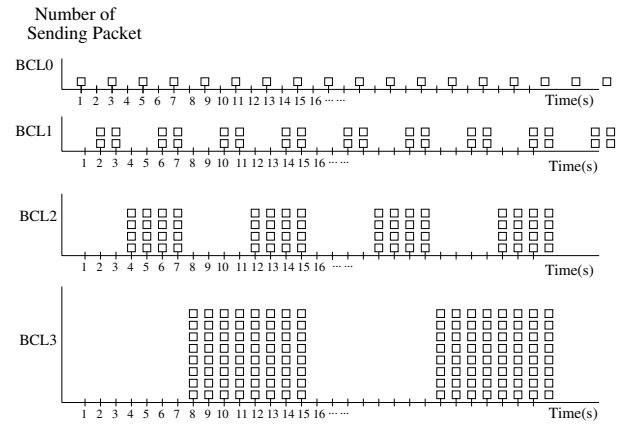


Fig. 2. Binary counting layers targeted for an RTT of 1 second

A. Introducing Binary Counting Layers

The key to our additive increase methods are binary counting layers, so named because the rates on the layers mimic aspects of counting in binary.

- *Binary Counting Layers (BCL)*: The rate transmitted on $BCL_i(x)$ is an on/off function with a sending rate of 2^i packets during each on time, and where the duration of each on and off time is $x \cdot 2^i$.

In TCP, the rate of additive increase is a function of the round-trip time: the window opens by one additional packet per RTT. The set of $BCL(x)$ layers provide the same functionality as TCP's additive increase with a measured RTT of x seconds. BCLs accommodate asynchronous join attempts for different receivers in the multicast session, and accommodate receivers with different target rates for the join attempt.

All layers are initially synchronized at time zero, which corresponds the beginning of an off time for all layers. Figure 2 shows how each Binary Counting Layer is organized, assuming a 1 second RTT, which we use throughout this discussion for simplicity.

To achieve additive increase starting at time zero, the receiver simply subscribes to BCL_i at $2^i \cdot \text{RTT}$ seconds. In Figure 2, where the RTT is 1 second, the receiver subscribes to BCL_0 , BCL_1 , BCL_2 , and BCL_3 at 1s, 2s, 4s, and 8s respectively. Once the receiver subscribes up through BCL_i , the number of receiving packets per RTT has increased by $2^{i+1} - 1$ with only i IGMP joins and no additional IGMP leaves. Avoidance of IGMP leaves is crucial, since in current versions of IGMP, it often takes a number of seconds before the leaves actually take effect; moreover, other extant methods for additive increase require use of IGMP leaves.

Previous work has defined *join* and *leave complexity*, i.e. the number of IGMP joins and leaves per operation, to be useful performance metrics for layered multicast [5]. For SMCC, the notion of operation does not map cleanly onto the additive increase process, so we will consider the complexity of N successive additive increases. From the description above, it is clear that this process requires $\log N$ joins (and no leaves). In

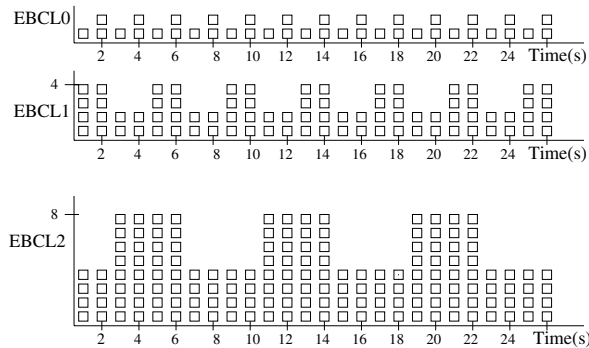


Fig. 3. Extended Binary Counting Layers

other approaches to additive increase, such as [6], the receiver periodically increases its rate by a constant amount c using a constant number of operations (typically 1 join and 2 leaves). Thus the complexity of N successive additive increases in these schemes is N/c , i.e. linear in N .

B. Extended Binary Counting Layers

One limitation of the basic binary counting layer scheme is that the receiver has to wait until certain specific times to join the BCLs. Suppose the receiver wants to increase its rate from 1 to 14 packets per RTT in Figure 2. If the receiver wants to join BCLs at 5 seconds, it has to wait until the next cycle (time 17) to initiate additive increase. One solution is to allow receivers to jump-start their additive increase with an initial set of joins (i.e. an immediate increase of 5 packets per RTT in the example above). However, this can induce sudden rate increase, and in the worst case, reduces to a naive join attempt. An alternative is the following improvement.

- *Extended Binary Counting Layers*: The rate transmitted on $EBCL_i(x)$ is a cyclic two step function. The number of sending packet during RTT x seconds is 2^i and 2^{i+1} in the low step and in the high step respectively.

Figure 3 shows the transmission rate of each binary counting layer targeted for an RTT of 1 second. The receiver subscribes to $EBCL_i$ at $(2^{i+1} - 1) \cdot RTT$ seconds to perform the additive increase. In Figure 3, the receiver subscribes to $EBCL_0$, $EBCL_1$, $EBCL_2$, and $EBCL_3$ at 1s, 3s, 7s, and 15s respectively to get the additive increase up to 30 packets per RTT. Now consider the waiting time if the receiver misses the join time. If the receiver has to start the increase from 1 packet to $2^i - 1$ packets, a new cycle for that increase starts at $2^i \cdot RTT$ seconds after the previous cycle starts in the basic BCL. However, in the extended BCL the new cycle for that increase starts at $(2^{i-1} + 1) \cdot RTT$ seconds after the previous cycle starts. For example, for the increase from 1 to 30 packets, the new cycle starts at 32 seconds and 17 seconds in the basic BCL and in the EBCL respectively.

So far we have accommodated receivers with a specific RTT (1 sec. in our example). In practice, receivers may have widely varying RTTs, and it is desirable to simulate TCP behavior of one packet per RTT additive increase for each receiver.

Extended BCLs can achieve this. The full version of this paper [15] shows how EBCLs can be organized to simultaneously accommodate receivers with various RTTs which are powers of two.

C. Cost of additional BCLs for join attempt

One cost of additional layers to facilitate additive increase is that they consume additional bandwidth beyond what is used by the normal cumulative layers. To measure this cost, we use the measure of dilation, defined in [5] and recapitulated here.

Definition 1: For a layering scheme which supports reception rates in the range $[1, R]$, and for a given link l in a multicast tree, let $M_l \leq R$ be the maximum reception rate of the set of receivers downstream of l and let D_l be the bandwidth demanded in aggregate by receivers downstream of l . The *dilation* of link l is then defined to be D_l/M_l . Similarly, the dilation imposed by a multicast session on tree T is taken to be $\max_{l \in T}(D_l/M_l)$.

Lemma 1: The worst case dilation of SMCC with single set of BCLs is 1.75.

Proof: Let us suppose the highest layer subscribed to by any downstream receiver is the j th layer. The maximum rate induced by the join attempt of a receiver k is $B_j - B_{j-2}$ when the following case holds: 1) the cumulative sending rate up through L_j is the maximum rate B_i , and 2) the cumulative sending rate up through L_{j-1} is slightly higher than the minimum rate B_{j-2} .

When an active receiver k in L_{j-1} has a calculated rate that is in the range of L_j , it performs a join attempt, which lasts until the total reception rate is equal to the next layer's cumulative sending rate B_j . Therefore, the maximum rate induced by the join attempt is $B_j - B_{j-2}$. The maximum reception rate of the set of receivers is B_j and the bandwidth demanded in aggregate by receivers is $B_j + B_j - B_{j-2}$. Therefore,

$$\text{dilation} = \frac{B_j + B_j - B_{j-2}}{B_j} = 1.75$$

Even though this worst-case dilation is not negligible, in practice it occurs only rarely (when a join attempt occurs across a bottleneck link); moreover, the average dilation during a join attempt is much smaller than this worst-case.

V. EXPERIMENTS

We have tested the behavior of SMCC using the ns simulator [16]. In most of the experiments we describe here, we use RED gateways, primarily as a source of randomness to remove simulation artifacts such as phase effects that may not be present in the real world. Use of RED vs. drop-tail gateways does not appear to materially affect the performance of our protocol. The RED gateways are set up in the following way: we set the queue size to twice the bandwidth-delay product of the link, set minthresh to 5% of the queue size and maxthresh to 50% of the queue size with the gentle setting turned on. Our TCP connections use the standard TCP Reno implementation provided with ns.

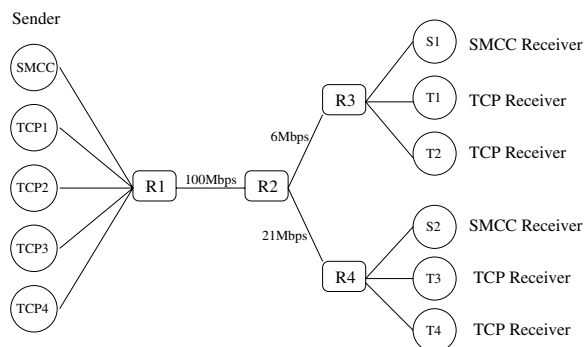


Fig. 4. Topology used to study TCP-fairness

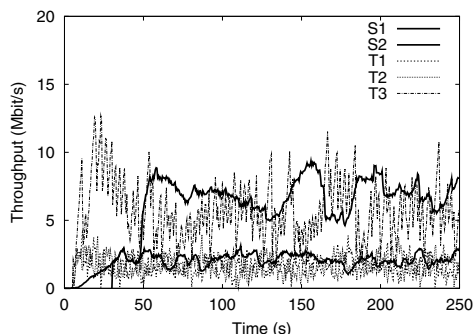
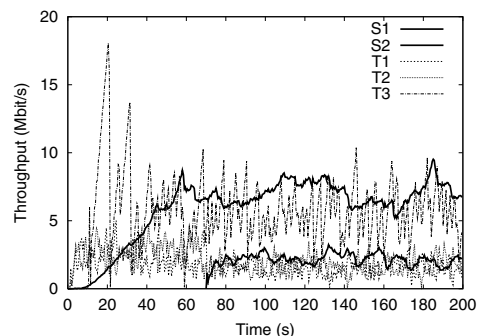


Fig. 5. Two SMCC receivers with TCP flows, $B_0 = 4\text{Mbps}$

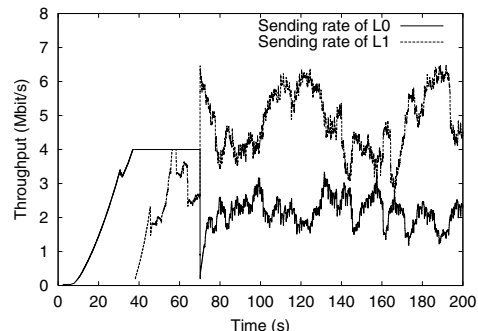
A. Preliminary Fairness Tests

Since the single rate TFMCC was well tested on the “dumbbell” topology [3], we set our initial topology to have multiple bottlenecks so that various SMCC receivers experience different network conditions. This initial topology is depicted in Figure 4. We set the propagation delay on each link is set to 8ms; each receiver therefore has a 64ms RTT in our simulations. Varying the delay on the links did not materially impact the performance of our protocol in the simulations we conducted. In this experiment and all of our other experiments we present, we used a damping factor of $\alpha = 1.2$. A full set of all the experiments we conducted as well as the ns source code are available online at <http://cs-people.bu.edu/guin/smcc.html>.

We consider a single SMCC session with two SMCC receivers and two parallel TCP flows sharing the same bottleneck link for each SMCC receiver. SMCC receiver S1 competes with 2 TCP connections on a 6Mbps link, giving a fair rate of 2 Mbps. S2 competes with 2 TCP flows on a 21Mbps link, for a fair rate of 7Mbps. We set B_0 to 4Mbps so that the sender’s maximum transmission rate on the base layer L_0 is 4Mbps. The throughput of each of the flows is plotted in Figure 5. S2 joins the base layer L_0 at 30.0 seconds, and it performs a join attempt at 47.6 seconds. After S2 subscribes to L_1 at 48.2 seconds, it shares fairly with the parallel TCP flows on the 21Mbps bottleneck link, while low-rate SMCC 1 shares fairly with 2 TCP flows on the 6Mbps link.



(a) Throughput of each receiver



(b) Sending rate of L_0 and L_1

Fig. 6. Late join of low-rate SMCC receiver. $B_0 = 4\text{Mbps}$

B. Late Join of Low-rate Receiver

In TFMCC, a late join by a low-rate receiver results in that low-rate receiver being selected as CLR, causing the sending rate of the entire session to be adjusted by its feedback. In SMCC, the late join of a low-rate receiver does not affect other receivers’ throughput on higher layers. Figure 6 (a) shows the throughput of SMCC receivers when the low-rate receiver, S1, joins late.

At the time S1 joins the session (70 seconds), the transmitted rate on the base layer is the maximum 4Mbps, while the rate on L_1 has been smoothly adjusting between 1 and 4Mbps to accommodate S2. The fair share for S2 behind the 6Mbps bottleneck link with two TCP competing flows is roughly 2Mbps, thus it immediately starts to experience a high loss rate. S1 is selected as C_0 within a second, and its feedback subsequently controls the transmission rate of L_0 . While the transmission rate of L_0 has changed from 4Mbps to S1’s feedback, the throughput of S2 is *not* adversely affected, since S2 is the CLR for L_1 , and the rate on L_1 instantaneously increases to compensate for the rate decrease on L_0 . Figure 6 demonstrates the discontinuities in the sending rates across L_0 and L_1 after time 70 seconds due to the late join of the low-rate receiver.

However, had there been other receivers subscribing only to the base layer, then the late join of a low-rate receiver clearly would affect other receivers at a same subscription level. The following rule is one of the keys to the scalability of our approach: degradation in the form of additional congestion

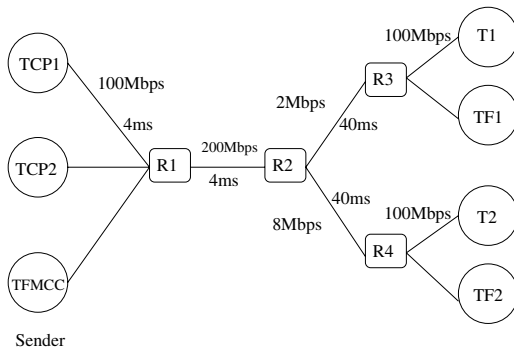


Fig. 7. Topology for calculated rate inaccuracy

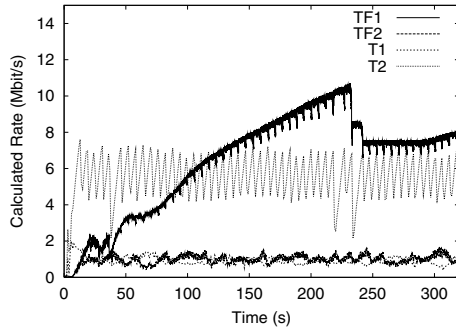


Fig. 8. Throughput of TF1 and TF2 over time. Competing TCP connections plotted in the background.

Receiver	Parameter	Time				
		50 s	100 s	150 s	200 s	250 s
TF1	RTT (second)	0.111	0.135	0.115	0.110	0.113
	Loss rate (%)	1.097	0.358	0.811	0.811	0.799
	Rate (Mbps)	0.761	1.175	0.880	0.919	0.902
TF2	RTT (second)	0.106	0.107	0.109	0.109	0.107
	Loss rate (%)	0.082	0.027	0.014	0.009	0.015
	Rate (Mbps)	3.220	5.583	7.671	9.366	7.442

TABLE II

CALCULATED TARGET RATE OF TFMCC RECEIVERS

along a path to a CLR will only impose throughput degradation to receivers at the same subscription level at that time. Rates received at other subscription levels are generally not impacted substantially.

C. Inaccuracy of Non-CLR Estimated Target Rate

Using TFMCC methods, a receiver which is not the CLR may not have sufficient information to correctly estimate its targeted rate. In particular, the loss rate measured by non-CLR receivers does not provide accurate information about the bottleneck bandwidth since the control equation was not modeled for this case, when the sender's transmission rate is independent of the receiver's packet loss events. The relevance of this point for SMCC is that a non-CLR receiver may not always be able to accurately assess whether it can safely join the next layer.

Figures 7 and 8 and Table II depict this scenario. In Figure 7, TFMCC receivers TF1 and TF2 are competing with two

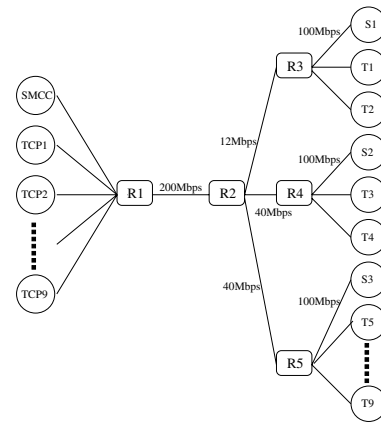


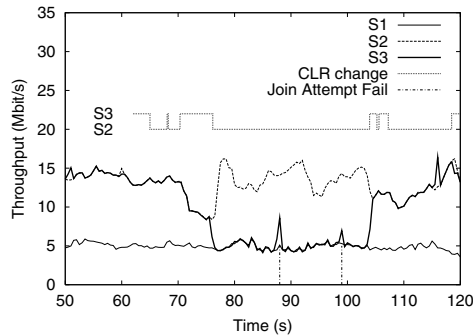
Fig. 9. Topology for assessing impact of dynamics of competing flow

TCP connections, T1 and T2, over a 2 Mbps bottleneck link and an 8 Mbps bottleneck link, respectively. TF1 and TF2 are not sharing the same bottleneck link, thus their losses are largely independent. Figure 8 shows each TCP flow's throughput and each TFMCC receiver's target rate calculated from the measured RTT and the loss rate. In the simulation, TF1 is quickly selected as C_0 and it fairly shares the 2 Mbps link with T1 throughout the simulation. Indeed, TF1's target rate over time, as depicted in Table II, is a reasonable approximation to its fair rate. In contrast, TF2's target rate, also depicted in Table II, is initially inaccurate (and badly underestimates the target rate) up through time 150 seconds. It then briefly overestimates its fair rate at time 200 seconds, and also overshoots its target subscription level, before converging around time 250 seconds. These estimation inaccuracies are another reason why we recommend and use conservative additive increase join attempts.

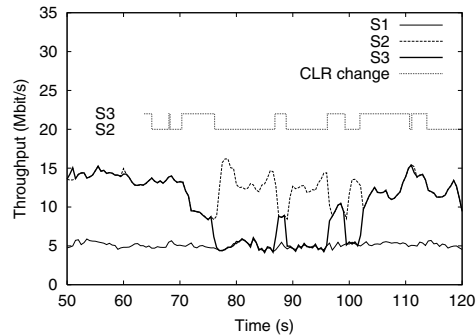
D. Responding to dynamics of competing traffic

We used a topology (Fig 9) to test the responsiveness to dynamic changes of local competing traffic, i.e. how increased traffic on local bottleneck links affects the receivers' throughput on different bottleneck links. As the competing traffic increases across a bottleneck, proportional fairness ensures that an SMCC receiver sharing the same bottleneck will get less throughput, and in the event that receiver is selected as CLR, the other receivers with the same subscription level also get less throughput even though they do not share the bottleneck with the CLR. However, the extent of the degradation is bounded by a penalty of at most a factor of 2 on all layers except for the base layer. Moreover, we will show that in practice, the degradation is typically much smaller than this worst-case bound.

In Figure 9, receiver S1 is competing with two TCP flows for a 12Mbps bottleneck link, while both S2 and S3 are competing with two different TCP flows for a different 40Mbps bottleneck link. We now set $B_0 = 8\text{Mbps}$ and all receivers have an RTT of 32ms. S2 and S3 do not share the same bottleneck link but their expected throughput is initially the



(a) With Additive Join Attempt



(b) Without Additive Join Attempt

Fig. 10. Impact of dynamic competing traffic, $B_0 = 8$ Mbps

same. Therefore, they have the same subscription level until new competing traffic starts.

Figure 10 (a) shows the throughput of each of the three SMCC receivers over time, as well as the CLR (either S2 or S3) on L_1 over time. Initially, the simulation starts with the three SMCC receivers and TCP flows 1 through 6. At 70 seconds, 3 additional TCP flows (T7, T8, T9) sharing the 40Mbps bottleneck enter the system. Therefore, S3's fair share drops from roughly 13Mbps to roughly 7Mbps. S3 is selected as C_1 at 70.3 seconds and the sending rate for L_1 steadily decreases, once it is controlled by its feedback. The receiver with the same subscription level, S2, suffers performance degradation as it gets the packets sent at the S3 feedback rate. But S2's receiving rate is adversely affected by the increase of traffic on the path to S3 only so long as S3 is C_1 . At time 75.7 seconds, S3 drops its highest layer, L_1 when its calculated rate drops to 7.74Mbps. S2 is elected as new CLR for L_1 at 76.2 seconds and its feedback controls the sending rate of L_1 , which then quickly rebounds. Meanwhile, L_0 continues to be limited by S1, who continues to have a lower fair share than S3, so S3 receives at a rate of approximately 5Mbps during this time.

Although S3's fair share is only 7Mbps, for reasons described in Section III-C, it cannot make a highly accurate assessment of its expected throughput while receiving at only 5Mbps, and these inaccurate estimates induce it to make join attempts to L_1 . S3 experiences two join attempts, both of which fail due to packet loss, between 70 seconds and 100 seconds. These two join attempts, marked by small spikes away from the S1 baseline, occur at 87.1 seconds and at 98.3 seconds. The little spikes around this time indicate these join attempt failures.

Finally, the three additional TCP flows leave at time 100 seconds. S3 performs a successful join attempt at 103.4 seconds and it reaches L_1 at 103.9 seconds, at which time it resumes sharing with S2.

Figure 10 (b) shows the identical simulation of each SMCC receiver but *without* the benefits of additive increase join attempts. Instead, in this simulation, the receiver naively joins an additional layer whenever the calculated rate is in the range

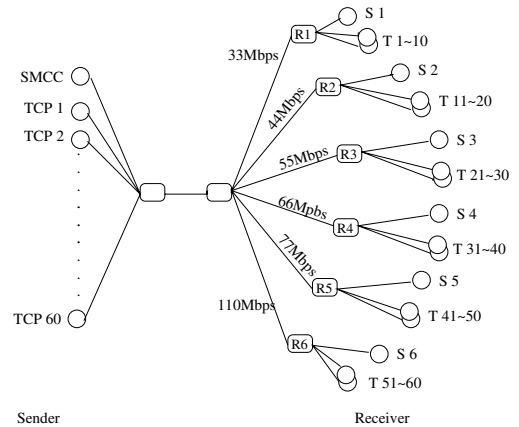


Fig. 11. Topology for Many Heterogeneous Receivers, $B_0 = 4$ Mbps

of the sending rate of the higher layer. S3 joins the next layer at 86.8 seconds and it becomes CLR for L_1 until 88.8 seconds. During this time, the sending rate of L_1 is dragged down to the rate of S3, impacting the reception rate of S2. After dropping back down, S3 joins L_2 at 96.1 seconds again and it is selected as C_2 until 99.3 seconds. Spurious joins such as these can cause significant performance degradation; an effect which is that much more severe when *multiple* receivers perform spurious joins, thereby constantly dragging down the rates on higher layers.

In contrast, with additive increase joins, even when a receiver initiates joins which are ultimately unsuccessful, it does not diminish the throughput received by other session participants during that time.

E. Fairness with heterogeneous receivers

Finally, we used a topology with multiple bottlenecks (Figure 11) to test the performance of SMCC with a set of heterogeneous receivers where the differences between the receivers' target rates is relatively small. We consider a single SMCC session with six SMCC receivers and ten parallel TCP flows sharing the same bottleneck link for each SMCC receiver, but each SMCC receiver is not behind the same bottleneck link. S1 competes with 10 TCP connections on a

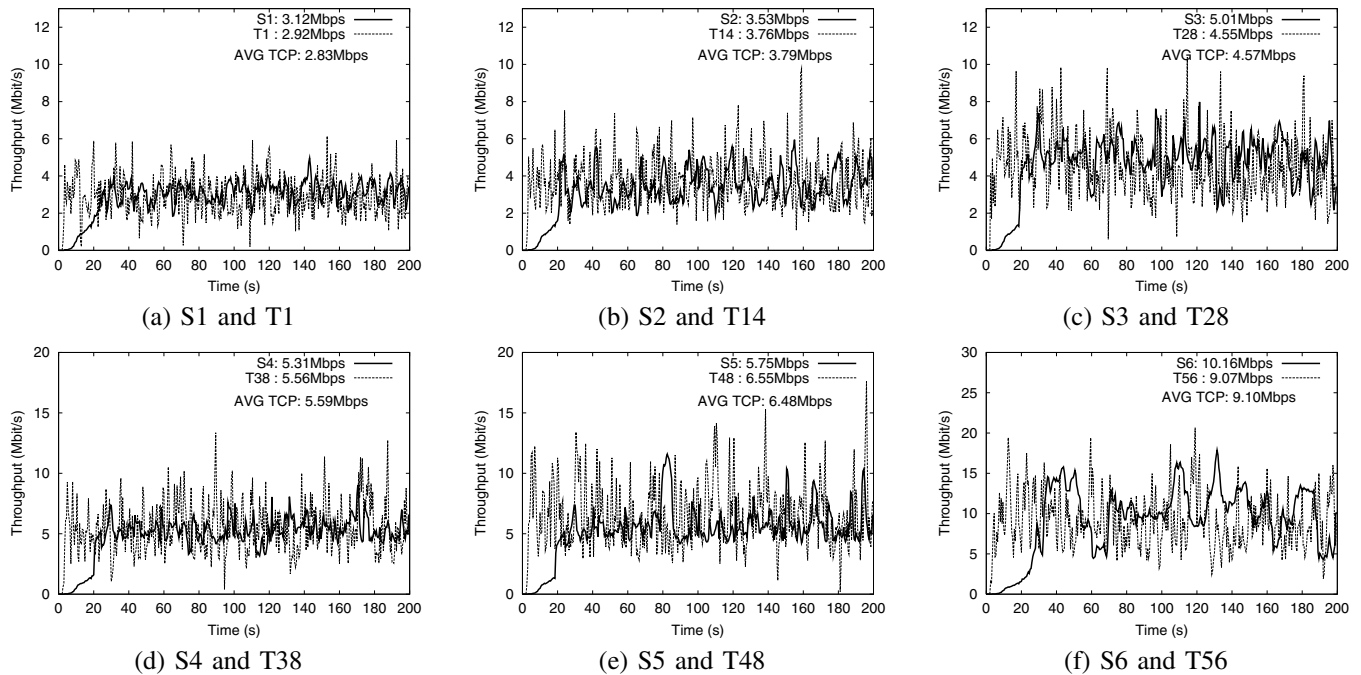


Fig. 12. Throughput of SMCC receivers, $B_0 = 4$ Mbps

33Mbps link, giving a fair rate of 3 Mbps and the fair rates of the other SMCC receivers (S2 to S6) are 4Mbps, 5Mbps, 6Mbps, 7Mbps, and 10Mbps respectively.

We plot the throughput of each SMCC flow and the throughput of one of the competing TCP flows in Figure 12. In each case, we chose the TCP connection whose mean rate was closest to the average of the ten competing flows as our representative. The throughput of each SMCC receiver fairly shares the bottleneck link with the parallel TCP flows even though lower-rate receivers are often present and drag down the rate on each level. In practice, non-CLR receivers tend to periodically join the next higher layer as their estimated throughput begins to deviate substantially from the CLR's target rate. The receiver S6 in panel (f) of Figure 12 is an example of relatively frequent subscription changes; note that its performance is still not as bursty as the competing TCP connection.

Next, consider Figure 13 (a) which plots the reception rate of S1 and S2. Like S6, S2 has relatively high subscription changes since its fair rate of 4Mbps is equal to B_0 . S1, with a fair rate of 3Mbps, is typically selected as the CLR on L_0 . Whenever S2 joins L_1 , it quickly becomes the CLR and may impact the throughput of receivers on that layer. The plot depicted in Figure 13 (b) shows this impact. For example, at time 67.08 seconds, S2 becomes C_1 and drags the cumulative rate T_1 down from 6.7Mbps to 5.0Mbps. At 99.6 seconds and 176.4 seconds, the sending rate is set from 7.1Mbps to 4.1Mbps and from 6.8Mbps to 4.6Mbps, respectively. There are other cases where the newly joined S2 becomes CLR on L_1 , but its degradation of rate is within 10%.

VI. CONCLUSION

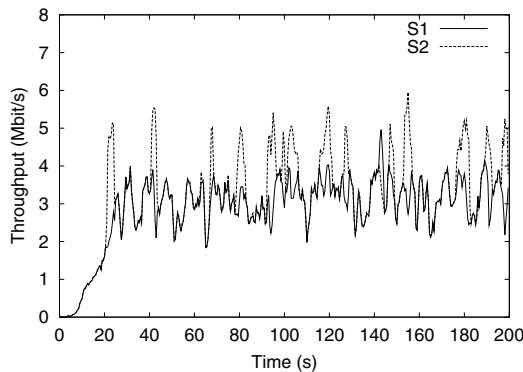
We have presented SMCC, a multirate equation-based multicast congestion control that leverages a proven single rate congestion control method (TFMCC) by orchestrating an ensemble of independently controlled single rate sessions. A compelling argument for this new methodology is its evident simplicity: unlike all other viable multiple rate congestion control protocols, ours requires only a small amount of carefully crafted new functionality. By maintaining appropriate invariants on the session rates of the individual TFMCC flows, specifying a clean mapping from reception rates to subscription levels and providing a non-disruptive method for additive increase join attempts, we build a sound multiple rate multicast congestion control scheme. A final advantage of our approach is its modular design; TFMCC could easily be replaced by an alternative, or an improved equation-based rate control mechanism.

ACKNOWLEDGMENTS

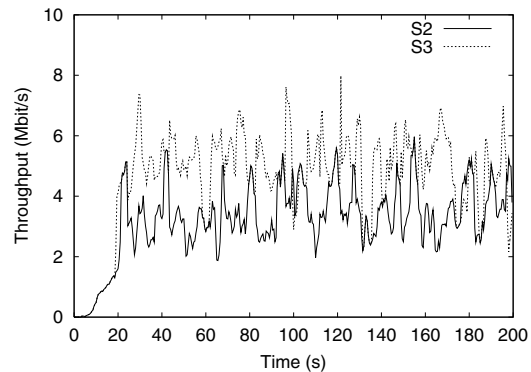
We thank the anonymous INFOCOM 2003 reviewers for their helpful comments. This work was supported in part by NSF grants ANI-9986397 and ANI-0093296.

REFERENCES

- [1] S. McCanne, V. Jacobson, and M. Vetterli, "Receiver-Driven Layered Multicast," in *Proc. of ACM SIGCOMM '96*, August 1996.
- [2] L. Rizzo, "pgmcc: A TCP-friendly single-rate multicast congestion control scheme," in *Proc. of ACM SIGCOMM '00*, 2000.
- [3] J. Widmer and M. Handley, "Extending equation-based congestion control to multicast applications," in *Proc. of ACM SIGCOMM '01*, 2001.



(a) Throughput of S1 and S2



(b) Throughput of S2 and S3

Fig. 13. Throughput Comparison of S1, S2, and S3

- [4] J. Byers, G. Horn, M. Luby, M. Mitzenmacher, and W. Shaver, "FLID-DL: Congestion Control for Layered Multicast," *IEEE J-SAC Special Issue on Network Support for Multicast Communication*, vol. 20(8), pp. 1558–1570, Oct. 2002. A preliminary version appeared in NGC '00.
- [5] J. Byers, M. Luby, and M. Mitzenmacher, "Fine-Grained Layered Multicast," in *Proc. of IEEE INFOCOM '01*, April 2001.
- [6] J. Byers and G. Kwon, "STAIR: Practical AIMD Multirate Multicast Congestion Control," in *Proc. of NGC '01*, 2001. Full version appears as BU-CS-TR-2001-018, Boston University, 2001.
- [7] A. Legout and E. Biersack, "PLM: Fast convergence for cumulative layered multicast transmission schemes," in *Proc. of ACM SIGMETRICS*, 2000.
- [8] L. Vicisano, L. Rizzo, and J. Crowcroft, "TCP-like Congestion Control for Layered Multicast Data Transfer," in *Proc. of IEEE INFOCOM '98*, April 1998.
- [9] M. Luby, V. Goyal, S. Skaria, and G. Horn, "Wave and Equation Based Rate Control Using Multicast Round Trip Time," in *Proc. of ACM SIGCOMM '02*, 2002.
- [10] B. J. Vickers, C. Albuquerque, and T. Suda, "Source-adaptive multilayered multicast algorithms for real-time video distribution," *IEEE/ACM Transactions on Networking*, vol. 8, no. 6, pp. 720–733, 2000.
- [11] J. Liu, B. Li, and Y. Zhang, "A Hybrid Adaptation Protocol for TCP-friendly Layered Multicast and Its Optimal Rate Allocation," in *Proc. of IEEE INFOCOM '02*, June 2002.
- [12] J. Padhye, V. Firoiu, D. Towsley, and J. Kurose, "Modeling TCP Reno Performance: A Simple Model and Its Empirical Validation," *IEEE/ACM Transactions on Networking*, vol. 8, no. 2, pp. 133–145, Apr. 2000.
- [13] J. Byers, M. Luby, and M. Mitzenmacher, "A Digital Fountain Approach to Asynchronous Reliable Multicast," *IEEE J-SAC Special Issue on Network Support for Multicast Communication*, vol. 20(8), pp. 1528–1540, Oct. 2002. A preliminary version appeared in ACM SIGCOMM '98.
- [14] S. Floyd, M. Handley, J. Padhye, and J. Widmer, "Equation-based congestion control for unicast applications," in *Proc. of ACM SIGCOMM '00*, 2000.
- [15] G. Kwon and J. Byers, "Smooth multirate multicast congestion control," Technical Report BU-CS-TR-2002-025, Boston University, 2002.
- [16] ns: UCB/LBNL/VINT Network Simulator (version 2). Available at <http://www-mash.cs.berkeley.edu/ns/ns.html>.