

FLID-DL: Congestion Control for Layered Multicast

John Byers ^{*} Michael Frumin [†] Gavin Horn [‡] Michael Luby [§]
Michael Mitzenmacher [¶] Alex Roetter ^{||} William Shaver ^{**}

ABSTRACT

We describe Fair Layered Increase/Decrease with Dynamic Layering (FLID-DL), a new multi-rate congestion control algorithm for layered multicast sessions. FLID-DL generalizes the receiver-driven layered congestion (RLC) control protocol introduced by Vicisano, Rizzo, and Crowcroft, ameliorating the problems associated with large IGMP leave latencies and abrupt rate increases. Like RLC, FLID-DL is a scalable, receiver-driven congestion control mechanism in which receivers add layers at sender-initiated synchronization points and leave layers when they experience congestion. FLID-DL congestion control coexists with TCP flows as well as other FLID-DL sessions and supports general rates on the different multicast layers. We demonstrate via simulations that our congestion control scheme exhibits better fairness properties and provides better throughput than previous methods.

A key contribution that enables FLID-DL and may be useful elsewhere is Dynamic Layering (DL), which mitigates the negative impact of long IGMP leave latencies and eliminates the need for probe intervals present in RLC. We use DL to respond to congestion much faster than IGMP leave operations, which have proven to be a bottleneck for prior work.

1. INTRODUCTION

One of the significant remaining hurdles to widespread adoption of IP multicast is the development of suitable congestion control algorithms. Ideally, one would hope for a

^{*}Boston University, Computer Science Department. This work was done while visiting Digital Fountain, Inc. Supported in part by NSF grant ANI-9986397. E-mail: byers@cs.bu.edu

[†]Stanford University and Digital Fountain, Inc.

[‡]Digital Fountain, Inc. E-mail: gavin@digitalfountain.com

[§]Digital Fountain, Inc. E-mail: luby@digitalfountain.com

[¶]Harvard University, Department of Engineering and Applied Sciences. This work was done while visiting Digital Fountain, Inc. Supported in part by the Alfred P. Sloan Research Fellowship and NSF CAREER grant CCR-9983832. E-mail: michaelm@eecs.harvard.edu

^{||}Stanford University and Digital Fountain, Inc.

^{**}Oregon Institute of Technology and Digital Fountain, Inc.

multicast analogue of TCP congestion control. Such a protocol would be an end-to-end congestion control mechanism that scales to large audience sizes, matches the functional relationship between throughput and packet loss rate at each receiver that TCP achieves, and provides responsiveness to changing network conditions on the order of a round-trip time, like TCP. Challenges include receiver heterogeneity, accurate modeling of TCP performance, and providing compatibility with other transport-level services such as reliability. In this paper, we provide a new multicast congestion control scheme that makes substantial strides toward a deployable solution.

Defining appropriate multicast congestion control algorithms which scale to large, heterogeneous audience sizes is essential for enabling multicast “killer apps” such as reliable content distribution to large audiences [3] and video streaming [10]. Multi-rate congestion control, as opposed to single-rate congestion control [15], is a de facto requirement for scaling to large audience sizes, to avoid the problem of establishing a single session rate which caters to the client demanding the smallest rate.

Standard approaches to multi-rate congestion control employ layered multicast [10, 16, 3] from a single source. Layered multicast organizes multiple multicast groups into logical layers. A host tunes its reception rate by subscribing to and unsubscribing from layers, i.e., by joining and leaving multicast groups. Different receivers may subscribe at different end-to-end subscription rates. A *cumulative* layered scheme has the additional property that all receivers must subscribe to and unsubscribe from layers in consecutive order. Several congestion control schemes for layered multicast sessions exist, but all have drawbacks.

Our congestion control scheme, which is based in part on the RLC protocol developed by Vicisano, Rizzo, and Crowcroft in [16], coexists with TCP, scales to large audience sizes, requires no changes to network routers or multicasting routing protocols, and faces no deployment hurdles (beyond those of deploying multicast in general). We call this scheme *FLID-DL*, which stands for “Fair Layered Increase/Decrease with Dynamic Layering”. FLID provides a generalization and simplification of RLC’s TCP-like congestion control mechanisms for an arbitrary organization of multicast layers; DL is a strategy which effectively avoids a major response bottleneck caused by IGMP leave latency.

The Asynchronous Layered Coding (ALC) approach within the Reliable Multicast Transport (RMT) working group of the IETF provides a reliable layered multicast solution for content distribution [8]. ALC currently lacks a congestion control protocol backed by consensus that is suitable for standardization within the IETF. We argue that FLID combined with DL provides a significant step towards a congestion control protocol for ALC.

The remainder of the paper is organized as follows. In Section 2 we discuss some of the basic issues associated with the design of multi-rate congestion control schemes and describe the RLC protocol in more detail. In Section 3, we define DL and demonstrate its capability to eliminate the performance penalty of slow IGMP leave operations. In Section 4, we define FLID and show how it provides improved TCP-friendly congestion control. In Section 5, we provide experimental evidence based on *ns* simulations to support our results, and we conclude with directions for future work in Section 6.

2. PREVIOUS WORK

The technique of congestion-controlled cumulative layered multicast was first proposed by McCanne, Jacobson and Vetterli [10] in the context of packet video transmission to large heterogeneous audiences. Their Receiver-Driven Layered Multicast (RLM) protocol achieves scalability by using a *receiver-driven* methodology, in which the hosts tune their subscription level by joining and leaving layers. They advocate an approach in which receivers periodically perform *join experiments* by subscribing to an additional layer and dropping a layer when they experience packet loss. There are several challenges that this approach introduces. First, one host's join experiments can introduce packet loss at other hosts behind the same bottleneck link, producing a potential source of unfairness or inefficiency. Second, standard approaches to cumulative layered multicast have exponentially increasing rates over the layers, which implies that the frequency of join experiments across the layers must be carefully designed to be friendly to TCP traffic and other sessions. Addressing these challenges motivated Vicisano, Rizzo, and Crowcroft to propose their Receiver-driven Layered Congestion Control (RLC) protocol [16].

2.1 RLC

RLC [16] was designed to provide a TCP-friendly multi-rate congestion control scheme which scales to large audience sizes, requires no modifications to routers or routing protocols, and does not require any coordination amongst receivers. For full scalability, a receiver-driven approach is required, as maintenance of per-receiver state at the source is infeasible and unscalable. But, uncoordinated join experiments by receivers pose substantial problems, as was observed in [10]. The authors of RLC cleverly avoid this problem by *synchronizing* join experiments. The source places synchronization points or increase signals into packets, where receivers can now only add a given layer after an appropriate increase signal for that layer. These increase signals are also *cumulative*, i.e., an increase signal j indicates that all receivers whose maximum subscription level is at most j can join a single additional layer. The use of cumulative increase signals solves the problem of synchronizing receivers behind a shared bottleneck, since when one receiver joins a layer that exceeds the bottleneck bandwidth, all other receivers behind that bottleneck will have also joined a layer. Then, since they will all experience packet loss, they will all drop back to their original rate prior to the join experiment.

In practice, care must be taken whenever a join experiment is performed, since by oversubscribing, a receiver can push the network into a state of congestion. To alleviate the congestion, the receiver must then unsubscribe from the layer by performing an IGMP leave operation, which can often incur substantial latency, leaving the network in

a congested state¹. Since oversubscription incurs a substantial cost, to minimize the likelihood of oversubscribing, the RLC source periodically injects a brief burst of packets on each layer prior to a synchronization point on that layer. The burst on layer i is designed to simulate the rate of layer $i + 1$, the idea being that those receivers which lose packets during the burst learn that adding layer $i + 1$ is unsafe, without incurring the cost of a join and leave operation. Unfortunately, if a receiver does *not* lose a packet in the burst, it still has no guarantee that adding the layer is safe, since the burst may be of insufficient length to induce packet loss (bursts recommended in [16] can be as brief as 8 packets). Thus a receiver is still prone to oversubscription. The complexity and lingering uncertainty associated with avoiding costly IGMP operations is one of the main problems with RLC which we address.

Another challenge addressed by RLC is the problem of appropriately orchestrating synchronization signals across the layers. The primary goal for RLC is to be fair to other instances of itself as well as to other congestion control algorithms such as TCP. As with most proposed layered multicast schemes, RLC requires that the rates on the layers must be exponentially spaced using a doubling scheme, i.e., the rates on the layers follow the pattern $1, 1, 2, 4, 8, \dots$. While dropping a layer with this scheme performs a TCP-like multiplicative decrease, adding a layer suddenly doubles the rate. Therefore, RLC cannot be TCP-like at a fine granularity, since it cannot perform fine-grained additive increase. However, it performs TCP-like additive increase at a coarser granularity by placing increase signals on layer i at a frequency of $\frac{1}{R_i}$, where R_i is the cumulative rate through layer i . When used in conjunction with a doubling scheme on the layer rates, the trajectory induced by this distribution of increase signals corresponds to linear increase over large time scales.

One issue which RLC does not adequately address is the dramatic fluctuations in network bandwidth consumption and the potential for rapid queue buildup that a doubling scheme can induce. We recommend the use of schemes which exhibit slower exponential growth in the layer rate, providing gentler transitions during join experiments.

2.2 TCP fairness

An increasingly widely accepted measure of TCP friendliness is to compare the steady state throughput of a flow along a path to the throughput that TCP would achieve along that path. For TCP traffic, a great deal of work has been done to determine the equation expressing the throughput as a function of the packet size, the packet loss rate, and the round-trip time along that path [9, 6, 7, 12]. It has been advocated that any new congestion control algorithm should exhibit the same steady state flow rate as suggested by the TCP equation [9, 1]; ensuring that the flows will then share available bandwidth fairly across a bottleneck link, since across the bottleneck link both streams experience the same packet loss rate. In addition, there is an emerging body of research that suggests designing congestion control algorithms to explicitly use the TCP equation [7, 9, 13].

In order for a congestion control scheme to be fair in this sense against TCP, the flow rate must have the same behavior as the TCP equation over large time scales. (Note

¹We review the root causes of slow IGMP leave operations in Section 3.1

that while this is a necessary condition, it may not be a sufficient condition, as issues of variance and the mechanics of rate changes come into play. We will address this issue only through simulation.) The TCP throughput rate R , in units of packets per second, can be approximated by the formula in [12]:

$$R = \frac{1}{RTT\sqrt{q}(\sqrt{2/3} + 6\sqrt{3/2}q(1 + 32q^2))} \quad (1)$$

where R is a function of the packet loss rate q , the TCP round trip time RTT , and the round trip time out value RTO , where we have set $RTO = 4RTT$ according to [7]. In both RLC and our work, since different multicast hosts have different end-to-end latencies from the server, and since the multicast analogue of RTT is not well defined, a target value of RTT , which we call the *aggressiveness factor*, is fixed in advance to generate a target rate R .

2.3 The use of a digital fountain for layered multicast

Since receivers join and leave layers over time in layered multicast, it is hard to control or predict precisely which packets they will receive. While this is not particularly problematic for appropriately encoded streaming content [10] that does not need to be transmitted reliably, scheduling packets across the layers in *reliable* multicast applications is a challenging problem. Recently there has been much work on integrating forward error correcting (FEC) codes into layered multicast as an end-to-end solution for scaling reliable multicast to audiences with heterogeneous download bandwidths. The benefit of using an encoded data stream is that it is no longer necessary to solve the difficult problem of delivering every single packet to every single host, thus admitting some flexibility into the scheduling of data packets onto layers over time [16, 3]. In our implementation, we use a digital fountain encoding [3] to generate an effectively unbounded number of different forward error correcting packets to be scheduled among the different layers. In this approach, as soon as the receiver receives enough distinct encoding packets, it can recover the original data, independent from the particulars of which layers it subscribed to over time.

3. DYNAMIC LAYERING (DL)

A significant limitation of current approaches to congestion controlled layered multicast is the timeliness of joining and leaving groups. With all existing schemes, rate increases can only be accomplished by joining one or more multicast groups; likewise rate reductions can only be accomplished by leaving one or more groups. Large join latencies are not especially problematic, but they introduce sluggish behavior which may cause the throughput of multicast sessions to underperform. Large leave latencies pose severe problems however, as they limit responsiveness to congestion and can create unfairness to sessions which react more quickly to congestion. In practice, IGMP leave latencies can in fact be very substantial, often on the order of several seconds.

3.1 IGMP leave latency

With the IGMP group membership protocol [4, 5], when a host wants to stop receiving content from a multicast group it sends a leave message to the last hop router. In general,

the last hop router does not track the number of hosts beyond the interface participating in a given multicast group, thus it must poll the hosts to determine whether any are still active before stopping the flow of packets. To provide reliability, the router typically polls up to three times before terminating flow to the group. In current implementations, each polling attempt can take from 1 to 3 seconds, for an aggregate leave latency of between 3 and 9 seconds. During this time, multicast traffic continues to flow through the last-hop router, even if no subscribers are present. Therefore, deployable congestion control algorithms for layered multicast must avoid relying on IGMP leaves to respond to congestion effectively, at least until faster IGMP leaves [14] are implemented.

3.2 DL overview

The key feature we use to achieve this goal is the use of *dynamic* layers, or layers whose rates change over time. Dynamic layers are distinguished from *static* layers over which the rate of packet transmission to the layer remains fixed for the duration of the session. The use of static layers necessitates explicit IGMP leaves to perform congestion control; use of carefully designed dynamic layers does not.

Our approach employs the following paradigm: the sender decreases the sending rate of each layer over time, and thus a receiver can decrease its reception rate quickly, simply by not joining any additional layers. In order for receivers to maintain a given reception rate they must periodically join layers at a moderate pace, as though they are on a treadmill. In order to increase their reception rate they must join additional layers beyond those needed to maintain a constant rate. With this general approach, slow leave operations do not affect the responsiveness to congestion.

3.3 Emulating cumulative layered schemes

We now demonstrate how to emulate any static cumulative layer scheme with a dynamic layer scheme. Suppose there are ℓ static layers with rates $r_0, \dots, r_{\ell-1}$ where zero is the index of the base layer, $\ell - 1$ is the index of the highest layer, and a receiver always subscribes to a cumulative set of layers starting from 0. Key parameters in designing the dynamic layer scheme are upper bounds on join and leave latencies.

We define J to be the worst case join latency and L to be the worst case leave latency. We assume that join latency is generally small (tens or hundreds of milliseconds) with small variance, while leave latencies can be much larger and much more highly variable. Now let s be an integer and T a real number satisfying $J < T < L < (s-1)T$. Our corresponding dynamic layer scheme uses $\ell + s$ dynamic layers. Each layer transmits at a fixed rate for a *time slot* of length T seconds.

Let $d_0, \dots, d_{\ell+s-1}$ be the $\ell + s$ dynamic layers and for convenience in describing the scheme, define $r_\ell = r_{\ell+1} = \dots = r_{\ell+s-1} = 0$. In the dynamic scheme, the transmission rate on layer d_j has rate $r_{(\ell+j-i) \bmod (\ell+s)}$ during time slot i . An equivalent interpretation is that during time slot i , the layer d_j carries the traffic corresponding to the static layer $\ell + j - i \bmod \ell + s$. Hence, each layer d_j has a period of $\ell + s$ time slots, where it begins transmitting at the highest rate $r_{\ell-1}$, drops sequentially through rates $r_{\ell-2}$ down to r_0 at time slot boundaries, and then transmits no packets for s time slots. The periodicity described above is necessary to efficiently utilize the multicast address space, as each new

layer corresponds to a separate multicast address.

Example 3.1 Consider a layering scheme where $\ell = 4$ and $r_i = i + 1$, for $i = 1, \dots, \ell - 1$. Let $J = 10$ ms and $L = 1.5$ seconds. If T is chosen to be 1 second, then $s = 3$. Figure 1 shows the rates on dynamic layers d_0 and d_1 for the first seven time slots.

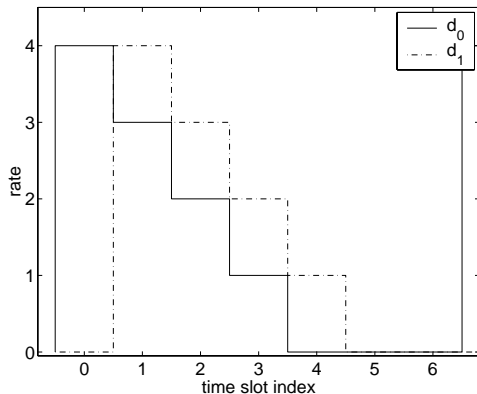


Figure 1: The rates on dynamic layers d_0 and d_1 for the first seven time slots.

In order to emulate the behavior of a static layer scheme, there are three operations which a receiver could perform and which must be emulated on the dynamic layers. For clarity, we assume that we have a receiver currently subscribing to dynamic layers d_k, \dots, d_{k+i} , i.e. emulating subscription to static layers $0, \dots, i$ in time slot k . (In the discussion of dynamic layering, all the indices are to be interpreted mod $\ell + s$.) The operations which a receiver can perform on the next time slot boundary are:

1. Emulate leaving static layer i . To do so, the receiver passively performs no action at the time slot boundary. The resulting aggregate rate will drop from $\sum_{j=0}^i r_j$ to $\sum_{j=0}^{i-1} r_j$.
2. Emulate retaining subscription to all current layers. To do so, the receiver joins dynamic layer d_{k+i+1} at the time slot boundary.
3. Emulate joining static layer $i + 1$: To do so, the receiver must join both dynamic layer d_{k+i+1} and dynamic layer d_{k+i+2} at the time slot boundary to move up to a rate of $\sum_{j=0}^{i+1} r_j$.

In all of the cases above, the receiver must also initiate a leave of dynamic layer d_k at the time slot boundary. We emphasize that this leave need not complete quickly, as layer d_k will transmit at a rate of zero for a substantial number of time slots. However, the leave must complete before the time slot at which reuse of layer d_k begins. Indeed, this explains why we use $\ell + s$ layers; we allow a layer to transmit at a zero rate for s consecutive time slots so that any leave message that occurred while the layer was transmitting has time to complete before the layer begins transmitting again at a non-zero rate. Note that when layer d_k is reused, it will start transmitting at the maximum possible rate.

The above also requires that the client know the current time slot (especially to correctly subscribe and leave dynamic layers). This can be accomplished in various ways, including using a separate multicast group for control information, or making the base layer static and embedding time slot information for the dynamic layers within the base layer.

To ensure that the reception rate of the receiver is smooth, a join can be scheduled just far enough in advance of the beginning of the time slot to ensure that packets start arriving from the joined dynamic layer just after the beginning of the next time slot. However, care must be taken that the join is not early enough to cause reception of packets from the joined dynamic group before the beginning of the time slot, as this may cause unnecessary congestion.

For the DL scheme, reactions (that is, joins and passive leaves) can occur at intervals of time T in a fast and predictable manner. The length of the time slot T is therefore a measure of the reactivity to network changes. T should be roughly the same as a small number of round trip times for TCP in order to be able to have approximately the same reaction times to changes in network conditions as TCP.

A receiver requires one leave and at most two joins to change its rate each T seconds. Smaller values of T make the system more responsive to loss; however, smaller values of T mean more join and leave requests per second. The overhead for these operations in terms of bandwidth and router utilization should be considered when designing the time slot and base layer bandwidth. Smaller values of T also mean more dynamic layers may be required. Larger values of T mean that on average there is a longer time interval between the time at which packet loss occurs and the time at which the receiver can respond to congestion.

3.4 Emulating non-cumulative layered schemes

Our results are not specific to cumulative layering approaches. An arbitrary static layering scheme consisting of ℓ layers can be emulated by a dynamic scheme with $s\ell$ layers, where s is as defined above. The idea is to use s dynamic layers to emulate each static layer; we cycle through those s dynamic layers so that only one is transmitting data during any time slot, and the transmission rate on the dynamic layer is the rate of the static layer (when it is not 0). Possible advantages of using non-cumulative layering schemes are discussed in [2].

4. FAIR LAYERED INCREASE/DECREASE

Fair Layered Increase/Decrease (FLID) is a protocol that is used to allow receivers to increase and decrease their reception rates based on congestion conditions so that the average throughput is similar to a TCP flow with a fixed RTT value, experiencing the same loss rate. FLID can either be implemented on top of DL, or on top of a static layering scheme (preferably in an environment where leave operations do not incur large latency).

FLID is akin to RLC [16] in many ways. The server places signals into packets that completely dictate the behavior of receivers with regards to joining and leaving layers. This property helps to coordinate the behavior of receivers behind bottleneck links. Like RLC, there is no feedback from receivers to the server, and different receivers may join different numbers of layers depending on the different network conditions on the paths between the server and the receivers.

These properties make FLID scalable to an unlimited number of receivers. In particular, receivers with slower bandwidth connections to the server do not slow down receivers with faster connections.

We introduce the FLID congestion control algorithm in three parts. First, we discuss the methods we employ to set the rates of the different layers. Next we present the actual rules for adding and dropping layers, and finally, we discuss how to set the increase signals at the server.

4.1 Rates on the different layers

FLID uses a cumulative layered scheme on ℓ layers. When a receiver subscribes to a set of layers $0, 1, \dots, i$, where $i \leq \ell - 1$, we call i the *subscription level*, or simply the level of the receiver. Generally the first layer is called the *base layer*, which in our case will be a static layer. To start receiving traffic for a multicast session, the host joins the base layer.

As before, r_i is the rate, in packets per second, for layer i and $R_i = \sum_{j=0}^i r_j$ is the cumulative rate, in packets per second, for layers 0 through i , i.e., the reception rate for a receiver with a subscription level i .

There are a variety of methods for choosing the different rates for the ℓ different layers. Some examples of schemes include:

- Equal scheme (E): The rates for all layers are equal, e.g., 3 packets per second. For example, with $\ell = 20$ layers and a 1 KB packet size, the subscription level can range from 24 Kbps up to 480 Kbps.
- Doubling scheme (D) [16]: The effect of adding another layer is to double the subscription level. The relative increases in the rates for the layers are in the sequence $1, 1, 2, 4, 8, 16, 32, 64, \dots$, i.e., $R_{i+1} = 2R_i$. For example, with a base layer rate of $R_0 = 3$ packets per second, $\ell = 20$ layers and a 1 KB packet size, the subscription level can range from 24 Kbps up to 12 Gbps.
- Multiplicative scheme (M): This is a generalization of the doubling scheme, where the rate for subscription level i is proportional to c^i for a fixed constant $c > 1$. The relative increases in the rates for the layers are in the sequence $1, c - 1, c^2 - c, c^3 - c^2, c^4 - c^3, \dots$, i.e. $R_i = c^i R_0$. For example, with a base layer rate of $R_0 = 3$ packets per second, $\ell = 20$ layers, a 1 KB packet size and $c = 1.3$, the subscription level can range from 24 Kbps up to 3.4 Mbps.

Other schemes are also possible. We focus on the multiplicative scheme for the remainder of the paper.

Two useful factors to consider in evaluating cumulative layered multicast schemes are the number of layers ℓ needed to span a given range of reception rates and the granularity with which a receiver can tune its rate within that range. In general, the tradeoff is that the larger the value of ℓ the more fine-grained rate changes can be and the smoother the reactions of the congestion control algorithm, but the more layers and hence multicast addresses that are needed for the transmission to achieve the same range of cumulative rates.

4.2 Increase and Decrease Rules

In FLID, the server partitions time into slots of duration T seconds each. All packets transmitted by the server in each time slot include the current *time slot index*. From the

receiver point of view, a new time slot starts when the first packet is received with a new time slot index.

Time slots are used to coordinate the activities of receivers. If during a time slot a receiver measures any packet loss, the receiver must decrease its subscription level by one at the end of the time slot. The receiver ignores all subsequent packet losses until the end of the time slot, i.e., the receiver only considers a single congestion event per time slot.

Time slots are also used to coordinate actions when receivers increase their subscription level. The server places an *increase signal* into each packet. The increase signal identifies a subscription level, which is fixed for all packets on all layers within a time slot. Receivers use increase signals to decide whether to increase their subscription level at the beginning of a time slot according to the following rule: if the current subscription level is i , the increase signal for the current time slot is j where $j \geq i$, and there is no packet loss during the current time slot, then increase the subscription level by 1 at the beginning of the next time slot. An increase signal of -1 is used to indicate that no receiver should increase its rate.

We use cumulative increase signals for the following reasons:

1. if a receiver behind a bottleneck link adds a new layer and does not experience congestion, then receivers sharing the same bottleneck link at lower subscription levels should also add a layer to fully exploit the available bandwidth;
2. if a receiver adds a new layer, causes congestion on a bottleneck link and drops back to the original subscription level, then receivers at lower subscription levels should have also added a new layer so that if they feel the same congestion, they also drop back to their original subscription level.

Of course there should be more increase signals for receivers at lower subscription levels than for those at higher subscription levels so that eventually all receivers behind the same bottleneck link will be at the same subscription level. The pattern of increase signals that the server uses is one of the more important considerations in a good design of FLID. This pattern is designed in conjunction with the rates of the different layers, and together they determine the fairness of FLID against other protocols such as TCP.

4.3 Setting the increase signals for the layers

To gain insight into how to set the increase signal values we consider a probabilistic pattern of increase signals. We then describe a deterministic pattern. Let $1 = p_{-1} > p_0 > p_1 > \dots > p_{\ell-1} = 0$. In each time slot, the increase signal is set to i with probability $p_i - p_{i+1}$ for $-1 \leq i \leq \ell - 1$. This means that a receiver with subscription level i will increase its level in each time slot with probability p_i . The value of p_{-1} must be set to zero, to prevent a receiver joined to all layers from ever attempting to join an additional layer that does not exist. For simplicity, we define $p_{-1} = 1$ which corresponds to a phantom layer that carries no traffic.

If we assume that each packet is lost independently with probability q , which is in fact the model used to derive equation (1), then for a given q and a set of p_i values we can calculate the average throughput \bar{R} as follows:

Define $P_j(k|i)$ to be the probability of a receiver moving from subscription level i to subscription level k at the beginning of time slot j . Since we can only increase or decrease the subscription level by a single layer at a time, $P_j(k|i) = 0$ for $|k - i| > 1$. In a time slot of length T , a receiver with subscription level i receives $T \cdot R_i$ packets, so we have

$$P_j(k|i) = \begin{cases} 1 - (1 - q)^{T \cdot R_i} & k = i - 1 \\ (1 - p_i)(1 - q)^{T \cdot R_i} & k = i \\ p_i(1 - q)^{T \cdot R_i} & k = i + 1 \\ 0 & \text{otherwise} \end{cases}$$

The transition probability $P_j(k|i)$ is independent of the time slot j so we write the transition probability as $P(k|i)$. (This is under the assumption that q does not vary over time.)

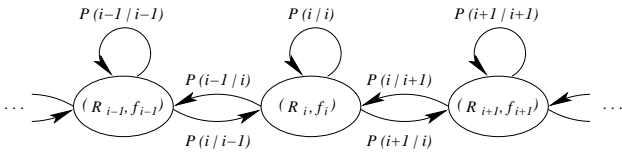


Figure 2: The Markov process showing the transitions for state (R_i, f_i) .

The receiver now behaves like an ℓ state Markov process a part of which is shown in Figure 2. We denote by f_i the steady-state fraction of the time the receiver will have a subscription level of i . At steady state, the average throughput \hat{R} is then

$$\hat{R} = \sum_{i=0}^{\ell-1} f_i R_i. \quad (2)$$

By the structure of the state diagram, at steady state the flow into each state on an edge has to be equal to the flow out of that state on that edge, i.e.,

$$P(i+1|i) \cdot f_i = P(i|i+1) \cdot f_{i+1}, \quad (3)$$

or

$$f_{i+1} = \frac{(1 - p_i)(1 - q)^{T \cdot R_i}}{1 - (1 - q)^{T \cdot R_{i+1}}} f_i \quad (4)$$

for $i = 0, \dots, \ell - 2$. Therefore, given a set of p_i and R_i values, it is simple to calculate \hat{R} using the recursion in (4) and equation (2).

To match the TCP equation, we have the following problem: given the target functional relationship between \hat{R} and q , find a set of strictly decreasing p_i values that approximate this function. In the appendix, we present an intuitive argument for setting the p_i values heuristically. We then develop a more sophisticated technique based on a hill-climbing algorithm.

Note the analysis above also immediately tells us the rate of join operations. Recall that each time slot requires a single leave operation and either 0, 1, or 2 join operations depending on whether the subscription level decreases, stays the same, or increases, respectively. From the above equations for $P(k|i)$, the expected number of joins per time slot when at subscription level i is $(1 + p_i)(1 - q)^{T \cdot R_i}$, and hence the steady state rate of join operations is $\sum_i f_i(1 + p_i)(1 - q)^{T \cdot R_i}$.

Example 4.1 Consider a base layer of $R_0 = 3$ packets per second, i.e., 24 Kbps for 1 KB packets, $\ell = 30$ and a multiplicative layering scheme with $c = 1.3$. For a time slot of $T = 0.5$ seconds and an aggressiveness factor of 0.1 seconds, if we set the p_i values according to the hill-climbing algorithm in the appendix, then we have the expected number of joins and leaves per second versus the packet loss rate q shown in Figure 3.

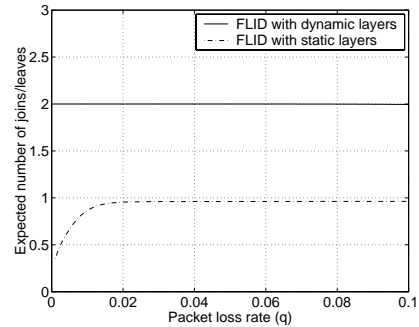


Figure 3: The expected number of joins and leaves per second versus the packet loss rate q for FLID with static layers and dynamic layers.

4.3.1 Use of a reverse binary counter

As mentioned previously, we use a deterministic scheme to set the increase signal in each time slot. This is because if we choose the increase signal in each time slot at random, then we occasionally increase the receive rate too rapidly when a receiver obtains several increase signals over a short time interval. Our deterministic scheme still generates an increase signal that allows a receiver to increase its subscription level from i to $i + 1$ about every $1/p_i$ time slots, but it has the additional property that the variance in the number of time slots between increase signals for each layer is minimal. We note that under the assumption that losses are independent, the steady-state probabilities of the Markov process will be the same for this deterministic scheme as for the probabilistic one, so we maintain throughput compatibility with TCP. Our scheme uses a reverse binary counter. Let

$$b = b_0 b_1 \dots b_{s-1}$$

be a number written in binary notation, i.e., b_i is a single bit and b is a concatenation of s bits. Let

$$\hat{b} = 0.b_{s-1} b_{s-2} \dots b_1 b_0$$

indicate b written backwards interpreted as a real number in the range $[0, 1]$.

At the beginning of a session, set b to be an arbitrary value, then, at each time slot, the algorithm for choosing the increase signal is to increment b by 1 and then find the largest layer i such that $p_i \geq \hat{b} \geq p_{i+1}$, and set the increase signal to i . We increment b modulo 2^s , so the counter is reset every 2^s time slots.

5. EXPERIMENTS

In order to show that FLID-DL is suitable for wide deployment on the Internet, we examine its behavior extensively

using *ns* [11]. We demonstrate two sets of experiments. In the first set, we consider how to set the various parameters of FLID-DL and show how FLID-DL addresses some of the shortcomings of RLC. In the second set, we examine how FLID-DL scales, measure its behavior to a set of heterogeneous clients, and demonstrate that it coexists fairly with different types of TCP traffic. We only have space for a summary of our findings here and refer the reader to www.digitalfountain.com/technology/library/flid for further details.

5.1 Setting the FLID-DL parameters

There are three parameters that we need to set for FLID-DL:

1. the rates on the different layers;
2. the aggressiveness factor; and
3. the duration of the time slots.

We implement FLID in our simulations with a multiplicative layering scheme to set the rates on the different layers. We use a base layer of 24 Kbps and a packet size of 1 KB, i.e., $R_0 = 3$ packets per second. For each independent FLID session we pick the number of layers ℓ large enough so that the highest reception rate possible for each individual session is greater than the capacity of the bottleneck link, up to a maximum of 30 layers.

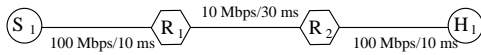


Figure 4: The network topology used to set the multiplicative factor c on the different layers. The node on the left is the server, the node on the right is the host, and the horizontal link in the middle is the bottleneck link.

To choose the rates on the different layers, we vary the value of the multiplicative factor c from 1.2 to 2.0. For each c value, we simulate a single FLID-DL session for 100 seconds and measure the throughput and number of packets lost over the final 50 seconds. Our topology consists of one server and one host connected by two DropTail routers, with a queue size of 128, and a 10 Mbps bottleneck between the two routers as shown in Figure 4. The time slot duration is 500 ms and the aggressiveness factor is chosen to be the same as the 120 ms round trip time of the topology.

Figures 5a and 5b show the throughput and number of packets lost in each 250 ms interval for $c = 1.3$ and $c = 2.0$ respectively. For $c = 1.3$, the bandwidth utilization is 85%, and an average of 47 packets are dropped per packet loss event. For $c = 2.0$, we had a bandwidth utilization of 59%, and an average of 17 packets dropped per packet loss event. However, over the 50 second interval, the latter experiment lost 2.5 times as many packets, as the large jumps in transmission rates cause more frequent packet loss.

Based on extensive simulations, of which this experiment is one representative, we chose $c = 1.3$ as our ideal multiplicative factor, since it is large enough to allow us to use few layers, but small enough to give us a sufficiently small granularity in the subscription level to avoid abrupt rate increases. It also gives a reasonably small packet loss rate when we exceed the bottleneck bandwidth.

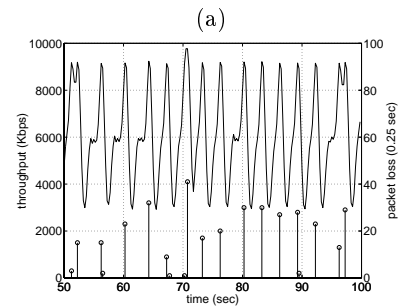
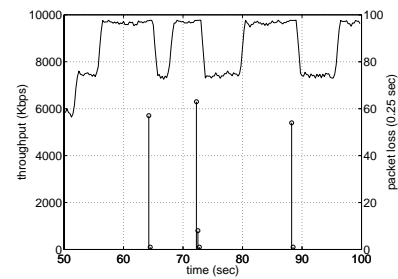


Figure 5: The throughput (solid line) and number of lost packets (circles) for a single FLID-DL session with (a) $c = 1.3$ and (b) $c = 2.0$.

For the rest of our experiments, we choose our aggressiveness factor to be 100 ms, so that FLID-DL competes fairly with TCP with an *RTT* of 100 ms. To set the ratio of the time slot duration to the aggressiveness factor, we use the same topology, except we now have random independent packet loss of 4% on the bottleneck link, i.e., we have no loss due to the queues overflowing. For a fixed aggressiveness factor, we vary the time slot duration and look at the throughput. For a given packet loss rate, doubling the time slot duration requires that we halve the reception rate in order to achieve the same probability of a packet being lost in each time slot. In order to achieve a reasonable reception rate; avoid large packet burst loss; and avoid a large number of join/leave operations per second, we choose the time slot duration to be 500 ms, i.e., 5 times the aggressiveness factor.

When simulating multiple TCP connections over a single bottleneck link in an event-driven simulator, synchronization problems can arise. To prevent this we start our connections at slightly different times and we add a small random delay to simulate processing overhead before each packet is sent out.

5.2 Static vs. dynamic layered schemes

Our next experiment compares the behavior of four independent FLID-DL sessions to four independent FLID-SL (static layer) sessions, when we have a random leave latency uniformly distributed between 2 and 4 seconds. Our topology consists of a simple double star network as shown in Figure 6. Both routers are DropTail with a queue size of 50 packets.

Figures 7a and 7b show the behavior of four FLID-DL sessions and four FLID-SL sessions with a leave latency of between 2 and 4 seconds respectively. The FLID-DL ses-

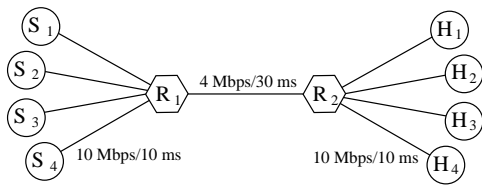


Figure 6: The network topology used to compare the static vs. dynamic layered schemes.

sions have a total bandwidth utilization of 90%, while for FLID-SL the bandwidth utilization is only 78%. In fact these results overstate the performance of FLID-SL, in that FLID-SL has reasonable bandwidth utilization only because in this simulation, hosts continue to accept packets received on groups which they have left but for which the router is still forwarding packets. The FLID-SL sessions combine to lose three times as many packets as their FLID-DL counterparts.

We have run the experiment with FLID-SL with various other leave latencies. When the leave latency is zero, we find that the behavior of FLID-SL and FLID-DL are very similar. The results are similar for RED routers with a queue size of 100 packets, a *maxthresh* of 50, a *minthresh* of 5 and gentle set to true.

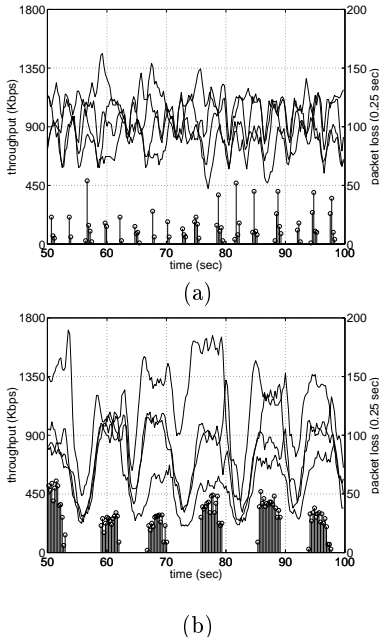


Figure 7: The throughput (solid line) and number of lost packets (circles) for (a) 4 FLID-DL sessions and (b) 4 FLID-SL sessions with a 2 to 4 second leave latency.

5.3 Coordination behind a bottleneck link

To demonstrate how hosts coordinate behind a common bottleneck link, we have 100 hosts subscribe to the same FLID-DL session at random times chosen uniformly between 0 to 5 seconds. We use the same topology as in Figure 6

except now with 100 hosts. Figure 8 shows that all 100 hosts converge to the same subscription level after 28 seconds.

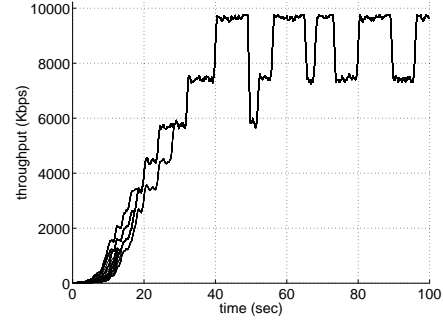


Figure 8: The coordination of 100 hosts behind a bottleneck link.

5.4 Random loss and heterogeneous delays

We test the scalability of FLID-DL in the presence of loss, where we have a number of hosts subscribed to a single FLID-DL session. We generate random loss at various points on each topology and measure the throughput downstream at the hosts. We find that the throughput for each FLID-DL host depends only on the loss rate experienced by that particular host. Experimental evidence indicates that FLID-DL scales well in the presence of random loss.

When we vary the delays for each host behind a common bottleneck, we find that the throughput at each host is proportional to both the aggressiveness factor and the random loss rate and is reasonably independent of the delay experienced by each host (except for the fact that the join at the highest layer takes longer to arrive). When the delay is greater than the time slot duration, FLID-DL has a subscription level that is higher than the actual reception rate. Since FLID-DL reacts to increase signals for its subscription level, as opposed to its actual reception rate, FLID-DL behaves less aggressively as the delay increases on the order of a time slot.

5.5 TCP fairness

For our final experiment, we compare how FLID-DL competes with TCP Reno and TCP SACK. We have n TCP and n FLID-DL streams share a common bottleneck of $0.5n$ Mbps. We vary n and calculate the throughput for the final 50 seconds of a 100 second simulation.

Figure 9a shows the relative throughput of FLID-DL and TCP for a DropTail queue of size $7n$ packets. Each point in the graphs represents the throughput of an individual stream. We also show the mean throughput for each type of stream. In this figure, which is representative of a large number of additional simulations, FLID-DL and TCP share the available bandwidth equitably.

Figure 9b shows the relative throughput of FLID-DL and TCP when we increase the DropTail queue size to $25n$ packets. In this scenario, FLID-DL is unfair to TCP and the average FLID-DL flow achieves a throughput 4 to 8 times larger than the average TCP flow. FLID-DL is slightly more fair in conjunction with TCP SACK.

FLID-DL is less fair to TCP as the queue size increases, since it does not adjust its reception rate in response to

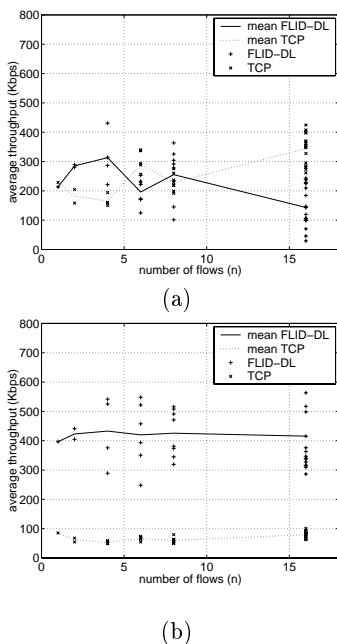


Figure 9: The fairness of FLID-DL and TCP Reno for a queue size of (a) $7n$ and (b) $25n$.

changes in the network round trip time. On the other hand, when the queue fills up, the delay increases and a TCP session slows down its transmission rate. Because of this, in an environment where varying queueing delays can account for a significant portion of the end-to-end latency, FLID-DL will not be fair to TCP.

6. CONCLUSION

We have demonstrated that the use of dynamic layering admits a simple solution to large IGMP leave latencies, without requiring changes to IGMP, routers, or other multicast routing protocols. We have also outlined the FLID scheme, which generalizes the RLC protocol but eliminates some of its complexity, such as probing. DL combined with FLID provides a significant step towards a complete and scalable receiver-driven congestion control algorithm for layered multicast. We hope that the introduction of FLID-DL will encourage additional work in multicast congestion control, and lead to quick standardization of the ALC approach in the RMT working group within the IETF [8].

Acknowledgements

The authors would like to thank Mark Handley and Lorenzo Vicisano for helpful suggestions in compiling a list of suitable experiments to test the FLID-DL protocol and the anonymous NGC reviewers for their suggestions for improving the paper.

7. REFERENCES

[1] B. Braden, D. Clark, J. Crowcroft, B. Davie, S. Deering, D. Estrin, S. Floyd, V. Jacobson, G. Minshall, C. Partridge, L. Peterson, K.K. Ramakrishnan, S. Shenker, J. Wroclawski, and

L. Zhang. Recommendations on Queue Management and Congestion Avoidance in the Internet. Technical Report IETF RFC 2309, April 1998.

[2] J. Byers, M. Luby, and M. Mitzenmacher. Fine-grained layered multicast. Unpublished manuscript submitted for publication.

[3] J. Byers, M. Luby, M. Mitzenmacher, and A. Rege. A digital fountain approach to reliable distribution of bulk data. In *Proc. ACM SIGCOMM*, Vancouver, Canada, September 1998.

[4] S. Deering. *Multicast Routing in a Datagram Internetwork*. PhD thesis, Stanford University, Stanford, CA, December 1991.

[5] W. Fenner. Internet group management protocol, version 2. Technical Report IETF RFC 2236, January 1997. Available at <ftp://ftp.isi.edu/in-notes/rfc2236.txt>.

[6] S. Floyd and K. Fall. Promoting the use of end-to-end congestion control in the Internet. *IEEE/ACM Trans. on Networking*, August 1999.

[7] S. Floyd, M. Handley, J. Padhye, and J. Widmer. Equation-based congestion control for unicast applications. In *Proc. ACM SIGCOMM*, 2000.

[8] M. Luby, J. Gemmell, L. Vicisano, L. Rizzo, J. Crowcroft, and B. Lueckenhoff. Asynchronous layered coding: A scalable reliable multicast protocol. Draft presented in the 47th IETF, March 2000.

[9] J. Mahdavi and S. Floyd. TCP-friendly unicast rate-based flow control. Note sent to end2end-interest mailing list, January 1997.

[10] S. McCanne, V. Jacobson, and M. Vetterli. Receiver-driven layered multicast. In *Proc. ACM SIGCOMM*, pages 117–130, Stanford, CA, August 1996.

[11] ns: Network simulator. Available at <http://www.isi.edu/nsnam/ns>.

[12] J. Padhye, V. Firoiu, D. Towsley, and J. Kurose. Modeling TCP throughput: a simple model and its empirical validation. In *Proc. ACM SIGCOMM*, Vancouver, Canada, September 1998.

[13] R. Rejaie, M. Handley, and D. Estrin. RAP: An end-to-end rate-based congestion control mechanism for realtime streams in the Internet. In *Proc. IEEE INFOCOM*, March 1999.

[14] L. Rizzo. Fast group management in IGMP. In *Hipparch Workshop*, pages 32–41, London, June 1998.

[15] L. Rizzo. pgmcc: a TCP-friendly single-rate multicast congestion control scheme. In *SIGCOMM 2000*, Stockholm, August 2000.

[16] L. Vicisano, L. Rizzo, and J. Crowcroft. TCP-like congestion control for layered multicast data transfer. In *Proc. IEEE INFOCOM*, San Francisco, CA, March 1998.

A. APPENDIX

A.1 Heuristic for setting the p_i values

We now provide a rough analysis of how to set the FLID parameters so that it behaves according to the TCP equation in the face of a fixed packet loss rate q . We model packet loss by a process where each packet is lost indepen-

dently with probability q , which is in fact the model used to derive the TCP equation (1).

Suppose the current subscription level is i , so that the aggregate rate is R_i . Since the probability that each packet is lost is q , on average there are roughly $t_d = 1/(qR_i)$ seconds between packet loss events. On the other hand, the increase signal when the subscription level is i occurs on average each $t_u = T/p_i$ seconds. If $t_d = t_u$, then the rate is as likely to go up as down for subscription level i . This occurs when $1/(qR_i) = T/p_i$, i.e., when

$$p_i = TqR_i. \quad (5)$$

We use this to set the values for the p_i 's since this approximately equates the probability of increasing and decreasing the rate in FLID, according to our Markov chain description in Section 4.3. We would roughly expect that if the steady state rate is R_i , then the rate is as likely to go up as it is to go down. While this is not precisely true since it depends on the relative values of the R_i 's, it is a good first approximation.

We now set the p_i values as follows:

1. Choose the aggressiveness factor for which we would like the FLID traffic to be fair.
2. For each subscription level R_i , we solve the TCP equation (1) for $q = q_i$ by setting R to R_i and RTT according to step 1.
3. We set $p_i = \min\{Tq_iR_i, 1\}$ for $i = 0, \dots, \ell - 1$.

When solving the TCP equation, the values of p_i may not be monotonically decreasing due to the influence of the RTT value on the TCP equation for large loss rates. To ensure that the FLID rules are followed, the p_i values for the small layers can be set to the maximum p_i value.

Example A.1 Consider a base layer of 24 Kbps, a packet size of 1 KB, i.e., $R_0 = 3$ packets per second, and a multiplicative layering scheme with $c = 1.3$ and $\ell = 30$. We set $T = 0.5$ seconds and the aggressiveness factor to be 0.1 seconds.

The graphs of the functions plotting average throughput (as derived from the TCP equation (1) and the FLID equation (2)) versus loss rate for TCP and FLID with these settings is given in Figure 10. As can be seen, the curves are quite close. We had similar results for other values of c as long as the ratio of the time slot duration to the aggressiveness factor was not too large.

A.2 Hill-Climbing Algorithm

The previous heuristic yields a set of p_i values that closely approximate the behavior of the TCP equation. However, because the analysis is clearly only approximate, we now describe an alternative approach that allows us to find a set of p_i values that provide a closer match, if this is desired, by applying a simple hill-climbing algorithm. Our approach here is more general in that it applies to any set of points describing the throughput for a given set of packet loss event rates.

A hill-climbing algorithm requires that we apply some metric to the p_i values as a measure of improvement. For this, we consider a specific range of loss rate probabilities $[q_a, q_b]$. We choose h linearly or logarithmically spaced points $q_a = q_1, q_2, \dots, q_h = q_b$, and evaluate the TCP equation (1) to determine the steady state transmission rate at these

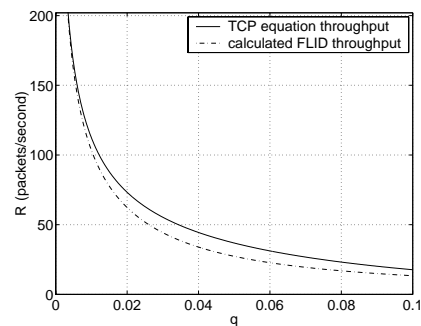


Figure 10: Comparison of the estimated average throughput as a function of the packet loss rate for FLID and TCP where the p_i values are obtained using the heuristic in Appendix A.1.

points for a given RTT . In fact we can evaluate any throughput equation for these loss rates. We denote the throughput versus packet loss rate curve as $R(q)$.

Given a set of p_i values we proceed as follows. For each q_j , $j = 1, 2, \dots, h$, we compute the average throughput $\hat{R}(q_j)$ at steady state by solving the Markov chain using the recursion in (4) and equation (2) in Section 4.3. We denote the resulting throughput versus packet loss rate curve as $\hat{R}(q)$. We then calculate the normalized mean square error ($NMSE$) between $\hat{R}(q)$ and $R(q)$ as follows,

$$NMSE(\hat{R}(q), R(q)) = \sum_{j=1}^h \frac{(\hat{R}(q_j) - R(q_j))^2}{R(q_j)^2}.$$

Of course we could also use the mean square error or another similar metric, such as the absolute difference, to measure the distance between the curve given by the cumulative layering scheme $\hat{R}(q)$ and the curve given by the TCP rate equation $R(q)$.

Our hill climbing algorithm is now as follows:

Hill Climbing Algorithm:
 Choose $[q_a, q_b]$ and RTT
 Calculate $R(q)$
 Initialize the p_i 's ($1 = p_{-1} > p_0 > \dots > p_{\ell-1} = 0$)
 do
 for $i = 0$ to $\ell - 2$
 Find the value of p_i ($p_{i-1} > p_i > p_{i+1}$)
 s.t. $\hat{R}(q)$ minimizes $NMSE(\hat{R}(q), R(q))$
 while (the p_i 's have not converged)

If we define p'_i to be the value of p_i before the execution of the for loop, then we say the p_i 's have converged if when the for loop completes we have

$$\sum_{i=0}^{\ell-2} (p_i - p'_i)^2 < \beta,$$

where β is a positive constant. In calculating the p_i 's, we chose $\beta = 0.00002$.

We note that in performing the hill-climbing algorithm, we maintain the restriction that the p_i values are decreasing.

Similarly, it may be desirable to ensure that the p_i values are somewhat separated, i.e. to ensure that $p_{i-1} - \alpha > p_i > p_{i+1} + \alpha$ for every p_i , where α is the minimum distance between consecutive p_i values. This ensures that it is possible for a layer to receive an increase signal without the layer above also receiving an increase signal. This prevents the unusual situation where subscribers up to layer i find they cannot increase their receive rate because every time they increase their receive rate, so do subscribers to the layer above, and the combination of increases causes significant loss. In calculating the p_i 's, we chose $\alpha = 0.001$.

A priori it is not clear how close we may come to the TCP curve. Experiments over a wide range of values have shown that the hill-climbing algorithm can very closely approximate the TCP curve with a cumulative layering scheme.

Example A.2 Consider a base layer of 24 Kbps, a packet size of 1 KB, i.e., $R_0 = 3$ packets per second, and a multiplicative layering scheme with $c = 1.3$ and $\ell = 30$. We set $T = 0.5$ seconds and the aggressiveness factor to be 0.1 seconds.

We use the hill climbing algorithm, where $h = 500$ and the q values linearly spaced in the range $[q_a, q_b] = [0.001, 0.1]$. The graphs of the functions plotting average throughput versus packet loss rate for TCP and FLID with these settings is given in Figure 11. As can be seen, the curves are very close.

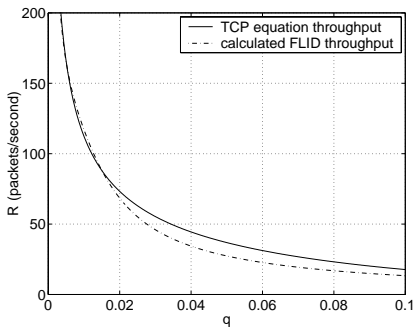


Figure 11: Comparison of the estimated average throughput as a function of the packet loss rate for FLID and TCP where the p_i values are obtained using the hill climbing algorithm in Appendix A.2.

More research is required to determine what p_i setting are appropriate in practice, and whether the $NMSE$ distance metric can be fine-tuned to take into account other considerations that may arise in the network setting.