

# A Self-Configuring RED Gateway

Wu-chang Feng†

Dilip D. Kandlur‡

Debanjan Saha‡

Kang G. Shin†

University of Michigan  
Ann Arbor, MI 48109

IBM T.J. Watson Research Center  
Yorktown Heights, NY 10598

{wuchang, kgshin}@eecs.umich.edu

{kandlur, debanjan}@watson.ibm.com

*Abstract*—The congestion control mechanisms used in TCP have been the focus of numerous studies and have undergone a number of enhancements. However, even with these enhancements, TCP connections still experience alarmingly high loss rates, especially during times of congestion. To alleviate this problem, the IETF is considering active queue management mechanisms, such as RED, for deployment in the network. In this paper, we first show that the effectiveness of RED depends, to a large extent, on the appropriate parameterization of the RED queue. We then show that there is no single set of RED parameters that work well under different congestion scenarios. In light of this observation, we propose and experiment with more adaptive RED gateways which self-parameterize themselves based on the traffic mix. Our results show that traffic cognizant parameterization of RED gateways can effectively reduce packet loss while maintaining high link utilizations under a range of network loads.

*Keywords*—Congestion control, Internet, TCP, RED, queue management.

## I. INTRODUCTION

Improving the congestion control mechanisms used in TCP has been an active area of research in the past few years. While a number of proposed TCP enhancements have made their way into implementations, TCP connections still experience high packet loss rates. Loss rates are especially high during periods of heavy congestion, when a large number of connections compete for scarce network bandwidth. Recent measurements have shown that the growing demand for network bandwidth has driven loss rates up across various links in the Internet [18].

One of the reasons for this high packet loss rate is the failure of the network to provide early congestion notification to sources. This has led the IETF to recommend the use of active queue management in Internet routers [1, 14]. There have been several active queue management algorithms which have been proposed including random drop [6, 13], early packet discard [21], and early random drop [8]. The most prominent and widely studied active queue management scheme is Random Early Detection (RED) [1, 5]. RED gateways alleviate many of the problems found in other active queue management algorithms in that they can prevent global synchronization, reduce packet loss rates, and minimize biases against bursty sources using a simple, low-overhead algorithm. Although RED outperforms traditional drop-tail queues, it is often difficult to parameterize RED queues to perform well under different congestion scenarios.

This paper demonstrates the importance of appropriately parameterizing RED queues by studying the performance of RED under a variety of configurations. The effectiveness of early detection is shown to be critically dependent upon the rate at which congestion notification is provided to the sources. When congestion notification is aggressive, packet loss can be avoided, but at the expense of low overall link utilization. On the other hand, when the rate of congestion notification is not sufficiently aggressive, link utilization remains high but the increased packet loss results in a drop in goodput and an increase in wasted resources. We present an on-line auto-configuration mechanism that determines suitable operating parameters for a RED queue depending on the traffic pattern. A key feature of this mechanism is that it does not require per-flow accounting, and retains the main advantages of RED. This is in contrast with other approaches which introduce additional complexity to the queue management scheme in return for additional functionality. Such schemes add either per-flow queueing [22] or per-flow accounting [11] in order to improve performance. Because of its simplicity and low overhead, the modifications proposed in this paper are particularly well-suited for managing traffic through large backbone routers where per-flow accounting and queueing may not be available.

Section II gives some background on how TCP congestion control and RED queue management work. Section III shows how the selection of RED parameters impacts performance across varying load. Section IV describes modifications which allow RED to self-parameterize itself in order to adapt to network load. In Section V, we present experimental results that demonstrate the efficacy of the proposed scheme. Finally, Section VI concludes the paper with a discussion of possible extensions.

## II. BACKGROUND

When a network is congested, a large number of connections compete for a share of scarce link bandwidth. Over the last decade, TCP congestion control has been used to effectively regulate the rates of individual connections sharing network links. TCP congestion control is window-based. The sender keeps a congestion window whose size limits the sending rate by limiting the number of unacknowledged packets that the sender can have outstanding in the network. When a connection starts

Parameter	Function
$min_{th}$	Threshold which the average queue length must exceed before any dropping or marking is done.
$max_{th}$	Threshold which the average queue length must exceed before all packets are dropped or marked.
$max_p$	Maximum marking probability. Determines how aggressively the queue drops or marks packets when congestion occurs.
$w_q$	Weight for updating average queue length. Large values increase the weight of more recent queue length data.
$f(t)$	Inter-packet drop probability function.

TABLE I  
RED PARAMETERS

up, it attempts to ramp up its transmission rate quickly by exponentially increasing its congestion window. This stage is called *slow-start* and allows the source to double its congestion window, and thus its transmission rate, every round-trip time. In order to prevent excessive losses due to an exponentially-increasing transmission rate, TCP senders typically employ what is known as the congestion-avoidance algorithm [7]. In this algorithm, TCP sources keep track of a threshold value ( $ssthresh$ ) which is a conservative approximation of the window size the network can support. When the window size exceeds this threshold, TCP enters the congestion avoidance phase. In this phase, the window is increased at a much slower rate of one segment per round-trip time. When the aggregate transmission rate of the active connections exceeds the network's capacity, queues build up and eventually packets are dropped. One way in which TCP detects a packet loss is through the receipt of duplicate acknowledgments from the receiver [9]. Upon receiving a small number of duplicate acknowledgments, TCP infers that a packet loss has occurred and immediately cuts its transmission rate in half by halving its congestion window. These mechanisms are called *fast retransmit* and *fast recovery*. When congestion is severe enough such that packet loss cannot be inferred in such a manner, TCP relies on a retransmission timeout mechanism to trigger subsequent retransmissions of lost packets. When a retransmission timer is triggered, TCP reduces its window size to one segment and retransmits the lost segment.

One problem with the TCP congestion control algorithm over current networks is that the TCP sources reduce their transmission rates only after detecting packet loss caused by queue overflow. This is a problem since considerable time may pass between the packet drop at the router and its detection at the source. In the meantime, a large number of packets may be dropped as the senders continue to transmit at a rate that the network cannot support. Active queue management has been proposed as a solution for preventing losses due to buffer overflow. The goal of active queue management is to detect incipient congestion *early* and convey congestion notification to the end-

hosts, allowing them to back off before queue overflow and packet loss occur. One of the more promising active queue management schemes being proposed by the IETF for deployment in the network is RED (Random Early Detection) [1, 4]. RED maintains an exponentially weighted moving average of the queue length which it uses to detect congestion. When the average queue length exceeds a minimum threshold ( $min_{th}$ ), packets are randomly dropped or marked with an explicit congestion notification (ECN) bit [4, 19, 20] with a calculated probability. When a TCP sender receives congestion notification in the form of an ECN mark, it halves its congestion window as it would if it had detected a packet loss. The probability that a packet arriving at the RED queue is either dropped or marked depends upon, among other things, the average queue length, the time elapsed since the last packet was dropped, and an initial probability parameter ( $max_p$ ). When the average queue length exceeds a maximum threshold ( $max_{th}$ ), all packets are dropped or marked. In our experiments, the aggressiveness of RED is varied by changing  $max_p$ , RED's initial drop probability. While it is possible to vary any of RED's parameters (as shown in Table I), we vary  $max_p$  since it directly impacts the aggressiveness of the early detection mechanism. One could potentially change  $f(t)$  and  $w_q$ , the inter-drop probability function and the queue weight respectively, however, doing so could potentially impact RED's fairness to bursty connections as well as its ability to avoid global synchronization effects.

### III. HOW EFFECTIVE IS RED?

The basic operating principle of all active queue management mechanisms is to provide early congestion notification to the sources in order to avoid packet loss due to buffer overflow. In heavily-congested networks, congestion notification must be signaled to a sufficiently large number of sources so that the offered load is reduced considerably. However, excessive early notification can lead to underutilization of network resources. Consequently, the success of an active queue management mechanism depends critically on how effectively it can strike a balance between reducing packet loss and preventing underutilization of the network by appropriately adjusting the rate of congestion notification.

For example, consider a bottleneck link of capacity  $10\text{Mbs}$  which is equally shared amongst several connections. When 100 TCP connections share this link, sending congestion notification to one connection reduces the offered load to  $9.95\text{Mbs}$ . On the other hand, when only 2 TCP connections share the link, sending congestion notification to one of them reduces the offered load to  $7.5\text{Mbs}$ . In general, when the bottleneck link is shared between  $N$  connections, a congestion notification signal to one connection reduces the offered load by a factor of  $(1 - \frac{1}{2N})$ . Hence, as  $N$  becomes large, the impact of individual congestion notifications decreases. In this situation, if the RED algorithm is not configured to be more aggressive, the RED queue can degenerate into a simple drop-tail queue. On the other hand, when  $N$  is small, the impact of individual con-

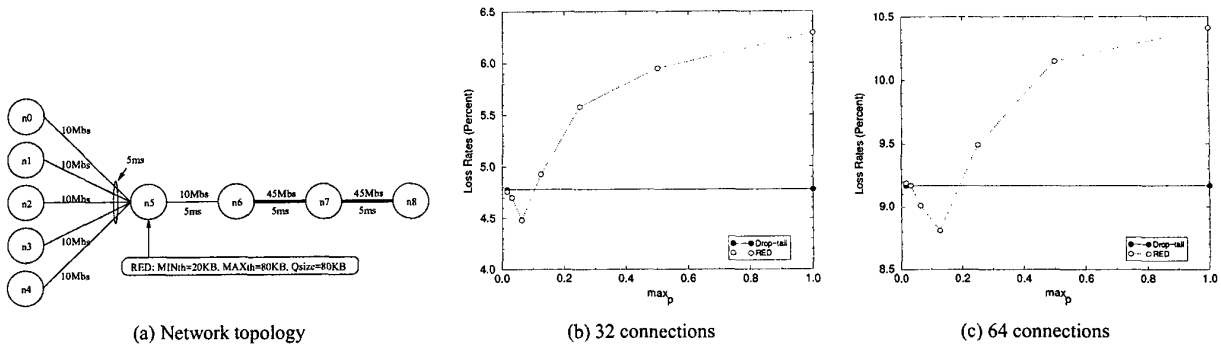


Fig. 1. Impact of early detection aggressiveness on RED

gestion notifications is large. In this condition, if the RED algorithm is aggressive, underutilization can occur when many sources back off their transmission rates in response to the observed congestion.

In the following sections, we analyze the effectiveness of early congestion notification used by RED. First, we consider the scenario where congestion notification is provided via packet drops. We show that RED is only marginally effective when used in conjunction with packet drop as a mechanism for congestion notification. Next, we study the performance of RED when used in conjunction with ECN enabled TCP connections. Although, RED is more effective when used with ECN enabled connections, packet losses are still significant. In either case, it is conspicuous that the congestion notification mechanism used by RED does not directly depend upon the number of connections multiplexed over the link.

#### A. RED with Packet Drop

In order to examine the impact of traffic load on early detection mechanisms used by RED, we performed a set of experiments using the ns simulator [15]. We varied both the aggressiveness of the early detection algorithm and the total number of connections multiplexed on the bottleneck link. Figure 1(a) shows the network topology used in these experiments. Note that each link is full-duplex, so acknowledgments flowing on the reverse path do not interfere with data packets flowing on the forward path. Each connection originates at one of the leftmost nodes ( $n_0, n_1, n_2, n_3, n_4$ ) and terminates at  $n_8$ , making the link between  $n_5$  and  $n_6$  the bottleneck.

Figure 1(b,c) shows the loss rates for 32 and 64 connections in a drop-tail and in a RED queue for a range of  $max_p$  values. In these experiments, packet drop is used to signal congestion to the source. This leads to an interesting optimization problem where the RED queue must pick a  $max_p$  value which minimizes the overall packet loss, including drops due to early detection and drops due to buffer overflow. When extremely large values of  $max_p$  are used, packet loss rates are dominated by drops due to early detection, while with extremely small values of  $max_p$  packet

loss is mostly due to queue overflow.

As the Figure 1 shows, RED has only a marginal impact on the packet loss rates observed. For small values of  $max_p$ , early detection is ineffective and the loss rates in the RED queue approach loss rates in the drop-tail queue. As  $max_p$  is increased, loss rates decrease slightly since the RED queue is able to send congestion notification back to the sources in time to prevent continual buffer overflow. Finally, as  $max_p$  becomes large, the RED queue causes a slight increase in packet loss rates over drop-tail queues. Note that the value of  $max_p$  that minimizes the loss rates is different in Figure 1(b) and Figure 1(c). As more connections are added, the optimal value of  $max_p$  increases.

The use of packet drops as a means for congestion notification fundamentally limits the effectiveness of active queue management. Steady state analysis of the TCP congestion avoidance algorithm [3, 10, 12, 15, 17] provides some insight as to why this is the case. Such analysis has shown that given random packet loss at constant probability  $p$ , the upper bound on the bandwidth of a TCP connection can be estimated as:

$$BW < \frac{MSS}{RTT} \frac{C}{\sqrt{p}} \quad (1)$$

where  $MSS$  is the segment size,  $RTT$  is the round-trip time, and  $C$  is a constant. Using this model, we can approximate packet loss rates over a single bottleneck link of  $L$  Mbs for a fixed number of connections  $N$ . In this situation, the bandwidth delivered to each individual connection is approximately the link bandwidth divided by the number of connections ( $L/N$ ). By substituting this into the previous equation and solving for  $p$ , we obtain

$$p < \left( \frac{N}{L} \frac{MSS * C}{RTT} \right)^2 \quad (2)$$

As the equation shows, when all connections use the TCP congestion avoidance algorithm, the upper bound on the packet loss rate increases quadratically with the number of connections present. Intuitively, this can be shown using an idealized example. Suppose we have two identical networks with bandwidth-delay products of  $64KB$

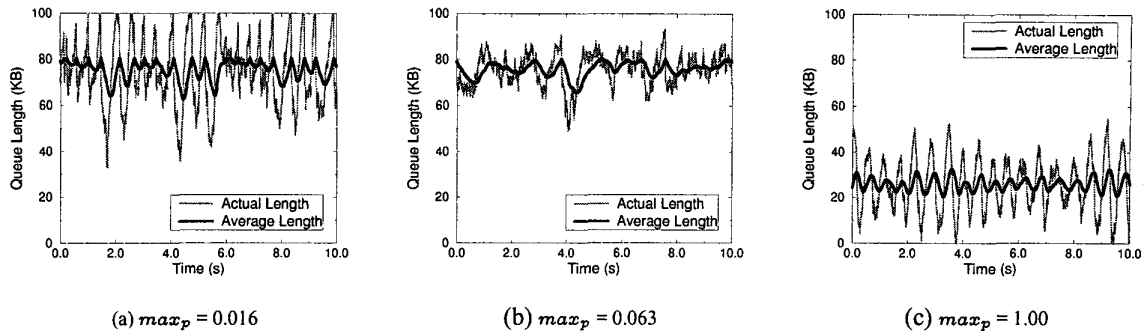


Fig. 2. Queue length plots of RED over  $max_p$

from a given pair of end points. In one network, we have 4 identical connections while in the other we have 8 identical connections going across the network. Given fair sharing amongst the connections, the congestion windows of each connection will approximately be the bandwidth-delay product divided by the number of connections present. For 4 connections, each connection will have congestion windows which oscillate near  $16KB$ . Assuming standard TCP windowing and a segment size of  $1KB$ , an individual connection in this network will typically build its congestion window up to around 16 packets, receive congestion notification in the form of a lost packet, back off its window to about 8 packets and then slowly build its congestion window back up at a rate of one packet per round-trip time. Given this behavior, the loss rate across all connections in this idealized model would then be approximately 4 packet drops every 8 round-trip times or 0.5 packets/round-trip time. Similarly, using the same idealized model, it can be shown that when 8 connections are present, losses occur at a rate of 2.0 packets/round-trip time, a quadratic increase.

It should be noted that the derivation of Equation 2 is based on idealized scenarios, the actual loss rates do not quite vary quadratically with the number of connections. From the drop-tail experiments in Figure 1, the observed loss rates show a dependence on the number of connections which is somewhere between linear and quadratic. This is partially due to the fact that connections can experience retransmission timeouts causing them to send less packets than the equation predicts. Still, since packet loss rates increase rapidly as networks become more congested, it is necessary to decouple packet loss from congestion notification through the use of ECN in order to prevent congestion collapse.

In lieu of explicit congestion notification, Equation 2 provides some insight on how to reduce loss rates. One way would be to eliminate the quadratic dependence on loss rate by making TCP's linear increase mechanism more conservative. For example, by making the increase linear in terms of the transmission rate, loss rates can be made to increase linearly with the number of connections. Without changing end-host congestion control, Equation 2 shows

that decreasing the segment size, increasing the bottleneck bandwidth, and increasing the round-trip times through the use of additional buffers [23] can slow the rate of congestion notification considerably and thus, decrease the amount of packet loss observed. Note that the dependence of packet loss rates on round-trip time explains why RED queues, which attempt to reduce network queuing delays, can potentially increase packet loss rates (as observed in Figure 1). To examine this particular effect, Figure 2 plots the queue lengths of the congested RED queue over several values of  $max_p$  in the experiment with 32 connections. When  $max_p$  is small as in Figure 2(a), early detection is ineffective and the behavior of the RED queue approaches that of a drop-tail queue. Figure 2(b) shows the queue length plot for the  $max_p$  value which minimizes packet loss (0.063). In this case, the RED algorithm is sufficiently aggressive to eliminate packet loss due to buffer overflow while keeping the average queue length as close to the size of the queue as possible. Finally, as  $max_p$  increases even further as shown in Figure 2(c), RED becomes more aggressive in its early detection leading to a drop in the average queue length. Consequently, the round-trip times seen by TCP connections also drop. Because packet loss rates increase with decreasing round-trip times (see Equation 2), this causes the loss rates in the RED queue to eventually exceed that of the drop-tail queue.

### B. RED with ECN

Next we examine the performance of RED using ECN-enabled TCP. By using ECN, all packet losses from the RED queue can be attributed to buffer overflow. We first focus on the effects of congestion notification in the desired operating range for RED where the queue length is between the  $min_{th}$  and  $max_{th}$  thresholds. This is achieved by setting  $max_{th}$  to the queue limit, which and causes packet loss to occur whenever early detection is not effective. In these experiments, the aggressiveness of RED algorithm is controlled using  $max_p$ .

Figure 3(a,b) shows a plot of the queue-length at the link between  $n_5$  to  $n_6$  when the number of connections active is 8 and 32 respectively. As Figure 3(a) shows, when only

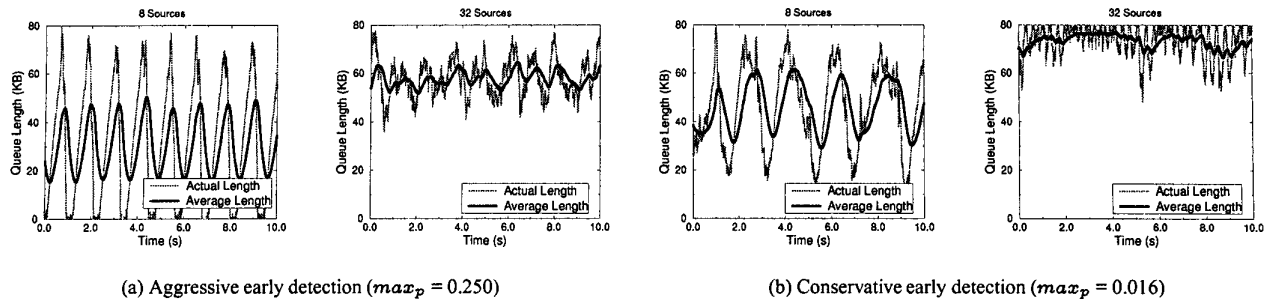


Fig. 3. Impact of early detection aggressiveness on queue length.

8 connections are active, aggressive early detection signals congestion notification to the sources at a rate which is too high, causing the offered load to be significantly smaller at times than the bottleneck link bandwidth. This causes periodic underutilization where the queue is empty and the bottleneck link has no packets to send. When the number of connections is increased to 32 the aggressive early detection performs as desired, sending congestion notification at a rate which can both avoid packet loss and achieve high link utilization.

Figure 3(b) repeats the same set of experiments, but with a more conservative early detection setting. In contrast to Figure 3(a), Figure 3(b) shows that by using less aggressive early detection, the RED queue can maintain high link utilization while avoiding packet loss when the number of connections is small. However, when the number of connections is increased to 32, Figure 3(b) shows that conservative early detection does not deliver congestion notification to a sufficient number of sources. Thus, the RED queue continually overflows causing it to behave more like a drop-tail queue. The figure also shows that the bottleneck queue never drains even though it is dropping a significant number of packets. This indicates that the TCP sources do not back off sufficiently in response to the congestion. In fact, the packets that are successfully delivered through the bottleneck queue trigger further increases in transmission rate. Consequently, the bottleneck queue remains close to fully occupied through the duration of the experiment.

To systematically evaluate the impact of  $max_p$ , we repeated the experiments across a range of traffic loads<sup>1</sup>. Figure 4(a) shows the loss rates observed for 4, 8, 32, and 64 connections. The figure also plots the loss rates when a drop-tail queue is used at the bottleneck link. As the drop-tail results show, loss rates at the bottleneck link increase proportionally to the number of connections using the link. The corresponding bottleneck link utilizations

are shown in Figure 4(b). The figures show that for small numbers of connections, loss rates remain low across all values of  $max_p$ , while only small values of  $max_p$  can keep the bottleneck link at full utilization. Thus, to optimize performance over a small number of connections, early detection must be made conservative. In contrast, for large numbers of connections, bottleneck link utilizations remain high across all values of  $max_p$  and only large values of  $max_p$  are able to prevent packet loss. In order to optimize performance in this case, early detection must be aggressive.

In the previous experiments, we set  $max_{th}$  equal to the queue size so that whenever the early detection algorithm fails, packet loss occurs. By setting  $max_{th}$  sufficiently below the queue size, the RED algorithm can avoid packet losses when early detection fails by deterministically marking every incoming packet. Figure 5 shows the queue-length plot using the same experiment as in Figure 3(b) with 32 connections, a larger bottleneck queue size and a fixed  $max_{th}$  of 80KB. When the queue size is 120KB, the queue-length plot shows that even with a fairly significant amount of additional buffer space, packet loss is not eliminated. It also shows that the combined effect of using ECN and packet drops for signaling congestion notification leads to periods of time where TCP is able to impact the transmission rates of the sources. This is in contrast to the behavior seen in Figure 3(b). In that experiment, whenever a connection is able to send a packet through the bottleneck link it increases its transmission rate even though the bottleneck queue is full. By setting  $max_{th}$  sufficiently low and using ECN, all connections receive congestion notification when the queue is full either in the form of ECN or packet loss. Thus, as Figure 5 shows, after a sustained period of ECN marking and packet loss, the sources back off sufficiently to allow the queue to drain. One of the problems with deterministic marking is that it often goes overboard in signaling congestion to the end-hosts. As the queue-length plots in Figure 5 show, periods of congestion are immediately followed by fairly long periods of underutilization where the queue is empty. Furthermore, the buffer space required to reduce loss is substantial. Figure 5(b) plots the queue-length using a queue size limit of 240KB. The figure shows that even though deterministic

<sup>1</sup>In each experiment, connections are started within the first 10 seconds of simulation. After 100 seconds, both the loss rates and the link utilization for the bottleneck link are recorded for 100 seconds. The loss rate is calculated as the number of packets dropped by the queue divided by the total number of packets which arrive at the queue. The link utilization is calculated as the total number of packets sent divided by the maximum number of packets the link could send.

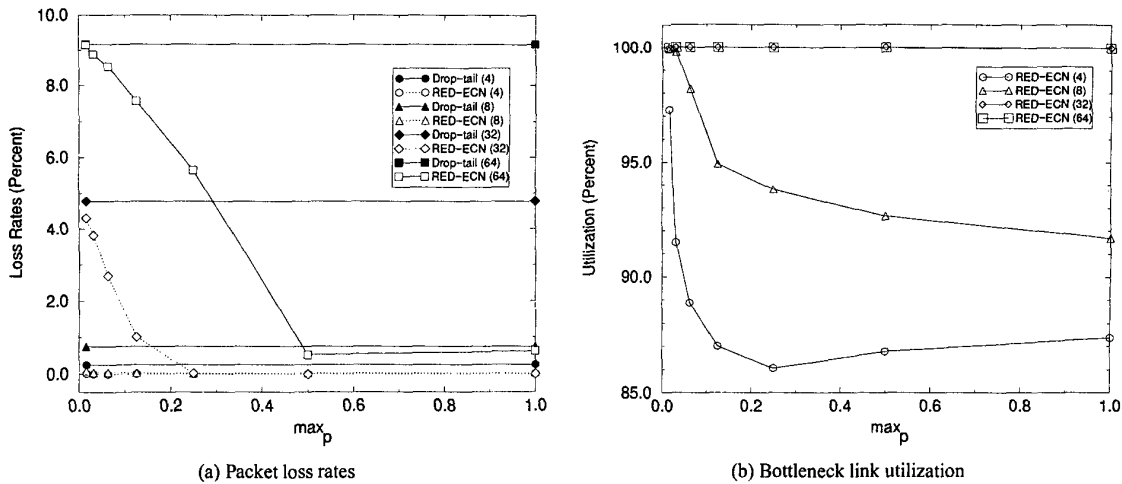


Fig. 4. Impact of early detection aggressiveness on RED-ECN

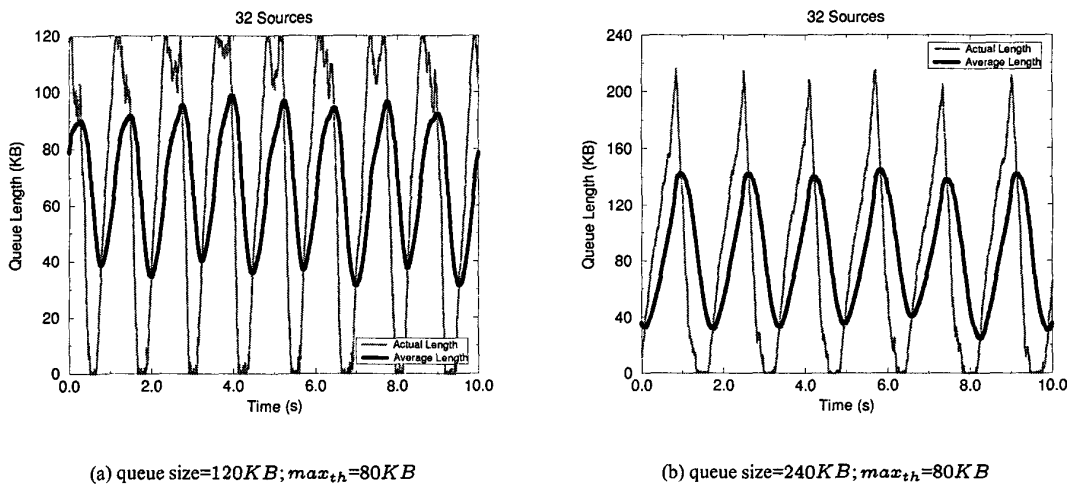


Fig. 5. Impact of  $max_{th}$  and queue size

marking is enforced at an average queue length of  $80KB$ , the actual queue length can more than double before the sources have a chance to back off.

#### IV. SELF-CONFIGURING RED

From the discussion in the previous sections, it is clear that adapting RED parameters based on traffic load and using explicit congestion notification can be beneficial to network performance. Figure 8 presents an on-line algorithm for adaptively changing the RED parameters according to the observed traffic. The idea behind this algorithm is to infer whether RED should become more or less aggressive by examining the variations in average queue length. If the average queue length oscillates around  $min_{th}$ , then the early detection mechanism is too aggressive. On the other

hand, if it oscillates around  $max_{th}$ , then the early detection mechanism is too conservative. Based on the observed queue length dynamics, the algorithm adjusts the value of  $max_p$ . In particular, it scales  $max_p$  by constant factors of  $\alpha$  and  $\beta$  depending on which threshold it crosses.

To explore the feasibility of the approach proposed in Figure 8, we ran another set of experiments using the same network as shown in Figure 1(a) with  $100KB$  queues. In this experiment, we vary the number of active connections between 8 and 32 over 40 second intervals. Figure 6 plots the queue-lengths for RED queues statically configured to be either aggressive or conservative. When aggressive early detection is used, as shown in Figure 6(a), the RED queue performs well whenever 32 connections are active. When only 8 connections are active though, the RED

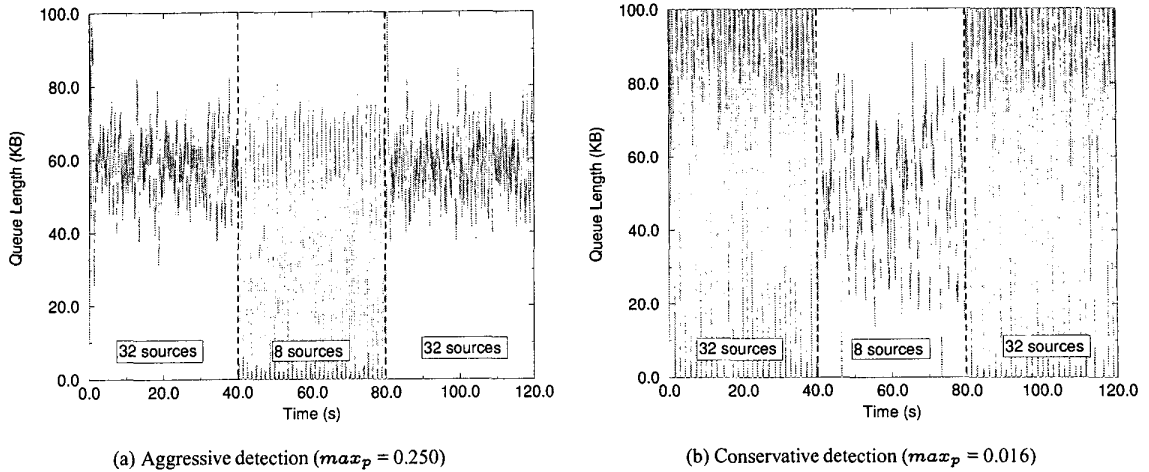


Fig. 6. Static random early detection

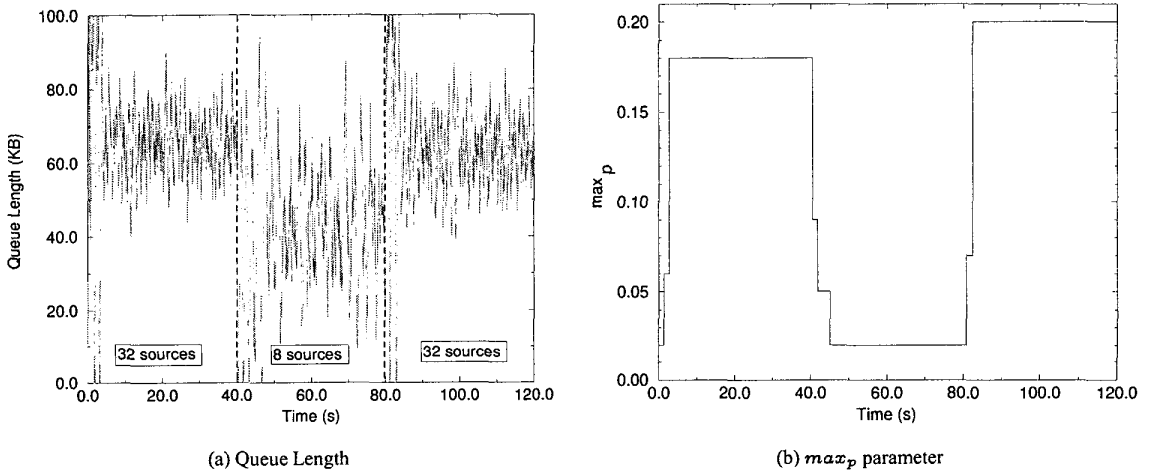


Fig. 7. Adaptive random early detection

```

Every  $Q(ave)$  Update:
  if ( $min_{th} < Q(ave) < max_{th}$ )
    status = Between;
  if ( $Q(ave) < min_{th}$  && status != Below)
    status = Below;
     $max_p = max_p / \alpha$ ;
  if ( $Q(ave) > max_{th}$  && status != Above)
    status = Above;
     $max_p = max_p * \beta$ ;

```

Fig. 8. Adaptive RED algorithm

queue is too aggressive in its congestion notification, thus causing periodic underutilization when the queue is empty.

When conservative early detection is used, as shown in Figure 6(b), the RED queue only performs well when 8 connections are active. When all 32 connections are active, the RED queue continually fluctuates between periods of sustained packet loss and ECN marking and subsequent periods of underutilization.

Figure 7(a) shows the queue-length plot of the same experiment using the adaptive RED algorithm with  $\alpha$  and  $\beta$  set to 3 and 2, respectively. We initially set  $max_p$  to 0.020 and allow the algorithm to adjust the parameter according to Figure 8. As the plot shows, after brief learning periods when the experiment starts and when the input traffic changes, the RED queue is able to adapt itself well. Figure 7(b) plots the  $max_p$  parameter as the RED queue adapts itself to the input traffic. As expected, its value adapts to reflect the number of active flows. When all 32 connections are active,  $max_p$  increases significantly, causing the

RED algorithm to become more aggressive. When only 8 connections are active,  $max_p$  decreases, causing the RED algorithm to become less aggressive.

## V. IMPLEMENTATION

We have implemented the adaptive RED modifications on a small testbed of PCs (Figure 9(a)) running FreeBSD 2.2.6 and ALTQ [2]. In this implementation, calls to the generic `IF_ENQUEUE` and `IF_DEQUEUE` macros from the `if_output` and `if_start` are changed in order to replace the FIFO queueing mechanism typically used in BSD UNIX with the adaptive RED queueing discipline. In addition, ECN is implemented using two bits of the type-of-service (ToS) field of the IP header [19]. When the adaptive RED algorithm decides that a packet must be dropped or marked, it examines one of the two bits to see if the flow is ECN-capable. If it is not ECN-capable, the packet is simply dropped. If the flow is ECN-capable, the other bit is set and used as a signal to the TCP receiver that congestion has occurred. The TCP receiver, upon receiving an ECN signal, modifies the TCP header of the return acknowledgment using a currently unused bit in the TCP flags field. A bit labeled “congestion experienced” (CE) is used by the TCP receiver to indicate the presence of congestion to the sender. Upon receipt of a TCP segment with CE bit set, the TCP sender invokes its congestion control mechanisms as if it had detected a packet loss.

Using this implementation, we ran several experiments on the testbed shown in Figure 9(a). Each network node and link is labeled with CPU model and link bandwidth, respectively. Note that all links are shared ethernet segments. Hence, the acknowledgments on the reverse path collide and interfere with data packets on the forward path. As the figure shows, FreeBSD-based routers using the adaptive RED queue management algorithm connect ethernet and fast ethernet segments.

We use `netperf` [16] to generate load on the system. `netperf` is a well-known, parameterizable tool for generating network load in order to evaluate the performance of both end-hosts and network elements. In the experiments shown here, a variable number of infinite `netperf` sessions is run from `fast` to `slow`, the endpoints of the network. The router queue on the congested interface on `router1` to the slow ethernet segment is sized at  $50KB$  using a  $min_{th}$  of  $10KB$  a  $max_{th}$  of  $40KB$ . A  $max_p$  value of 0.02 is used for the conservative early detection algorithm while a  $max_p$  value of 1.00 is used for the aggressive early detection algorithm. The  $max_p$  value of the adaptive algorithm is initially set at 0.02 and allowed to vary according to the algorithm in Figure 8. In order to ensure the adaptive RED modifications did not create a bottleneck in the router, we ran several experiments between `fast` and `router2` using the modifications on `router1` to forward packets between both hosts. In all of the experiments, the sustained throughput of `router1` was always above  $70Mbps$ . Thus, experiments run from `fast` to `slow` always bottleneck at the output interface to the slow ethernet segment on `router1`.

Figure 9 shows the throughput and packet loss rates at the bottleneck link across a range of workloads. The throughput measures the rate at which packets are forwarded through the congested interface. This rate slightly overestimates the end-to-end goodput of the `netperf` sessions since retransmissions are counted. The packet loss rate measures the ratio of the number of packets dropped at the queue and the total number of packets received at the queue. In each experiment, throughput and packet loss rates were measured over ten 5-second intervals and then averaged. As Figure 9(a) shows, both the conservative and adaptive early detection algorithms maintain high throughput levels across all workloads while the aggressive early detection algorithm achieves a lower throughput for smaller numbers of connections. Note that since the ethernet segments are shared, acknowledgments on the reverse path collide with data packets on the forward path, thus limiting throughput. Figure 9(b) shows the packet loss rates over the same workloads. As the figure shows, both the aggressive and the adaptive early detection algorithms maintain low packet loss rates across all workloads while the conservative early detection algorithm suffers from fairly large packet loss rates as the number of connections increase. As the number of connections increases, an interesting observation is that marking rate of the adaptive and aggressive early detection algorithms approaches 50%. Because aggregate TCP behavior becomes more aggressive as the number of connections increases, it becomes more and more difficult for the RED queue to maintain low loss rates. Fluctuations in queue lengths occur so abruptly that the RED algorithm oscillates between periods of sustained marking and packet loss and periods of minimal marking and link underutilization.

## VI. CONCLUSION AND FUTURE WORK

This paper has shown how adaptive active queue management can be used in conjunction with explicit congestion notification to effectively reduce packet loss in congested networks. In particular, an adaptive RED mechanism which is cognizant of the number of active connections can provide significant benefits in terms of decreasing packet loss and increasing network utilization. Preliminary experiments using a small testbed has shown the efficacy of the algorithm. Since early detection algorithms become more effective as round trip times increase, it is expected that the performance improvements will become even more significant as the network grows.

We are currently examining ways to methodically improve the adaptiveness of the RED algorithm. This paper presents one specific algorithm for tailoring RED parameters to the input traffic. There are many other potential alternatives for doing so. For example, the RED queue can actually keep track of the number of active connections and change its aggressiveness accordingly. Another mechanism would be to infer the number of connections present by the rate that the average queue length changes, then have the RED queue adapt its parameters accordingly. It may also be possible to adapt other RED parameters such



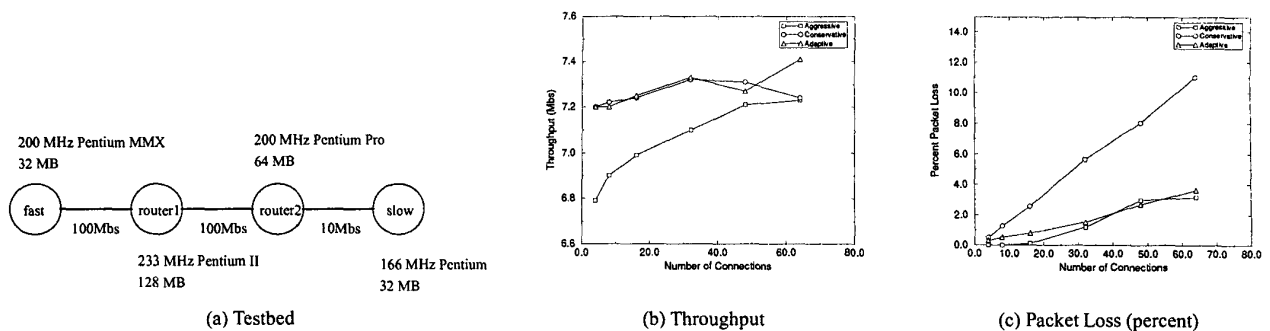


Fig. 9. Queue Management Performance

as inter-packet drop probabilities, RED threshold values, or queue length weights in order to optimize performance. For example, one could adapt the weight used to calculate the exponentially weighted average of the queue length based on the input traffic. When a large number of connections are present, queue lengths can fluctuate considerably in a short period of time. In such cases, an average which heavily weights more recent queue lengths is necessary in order to trigger congestion notification in time to prevent packet loss.

#### REFERENCES

- [1] R. Braden, D. Clark, J. Crowcroft, B. Davie, S. Deering, D. Estrin, S. Floyd, V. Jacobson, G. Minshall, C. Partridge, L. Peterson, K. Ramakrishnan, S. Shenker, J. Wroclawski, and L. Zhang. Recommendations on Queue Management and Congestion Avoidance in the Internet. *Internet Draft draft-irtf-e2e-queue-mgt-00.txt*, March 1997.
- [2] K. Cho. A Framework for Alternate Queuing: Towards Traffic Management by PC-UNIX Based Routers. *USENIX 1998 Annual Technical Conference*, June 1998.
- [3] S. Floyd. Connections with Multiple Congested Gateways in Packet-Switched Networks, Part 1: One-way Traffic. *Computer Communication Review*, 21(5), October 1991.
- [4] S. Floyd. TCP and Explicit Congestion Notification. *Computer Communication Review*, 24(5):10–23, October 1994.
- [5] S. Floyd and V. Jacobson. Random Early Detection Gateways for Congestion Avoidance. *ACM/IEEE Transactions on Networking*, 1(4):397–413, August 1993.
- [6] E. Hashem. Analysis of Random Drop for Gateway Congestion Control. *MIT Technical Report*, 1990.
- [7] V. Jacobson. Congestion Avoidance and Control. In *Proceedings of ACM SIGCOMM*, pages 314–329, August 1988.
- [8] V. Jacobson. Presentations to the IETF Performance and Congestion Control Working Group. August 1989.
- [9] V. Jacobson. Modified TCP Congestion Avoidance Algorithm. *end2end-interest mailing list (ftp://ftp.ee.lbl.gov/email/vanj.90apr30.txt)*, April 1990.
- [10] T. V. Lakshman and U. Madhow. The Performance of TCP/IP for Networks with High Bandwidth-Delay Products and Random Loss. *IFIP Transactions C-26, High Performance Networking*, pages 135–150, 1994.
- [11] D. Lin and R. Morris. Dynamics of Random Early Detection. In *Proc. of ACM SIGCOMM*, September 1997.
- [12] M. Mathis and J. Semke and J. Mahdavi and T. Ott. The Macroscopic Behavior of the TCP Congestion Avoidance Algorithm. *Computer Communication Review*, 27(1), July 1997.
- [13] A. Mankin. Random Drop Congestion Control. In *Proceedings ACM SIGCOMM (Special Issue Computer Communication Review)*, pages 1–7, September 1990. Published as Proc. ACM SIGCOMM '90; (Special Issue Computer Communication Review), volume 20, number 4.
- [14] A. Mankin and K. K. Ramakrishnan. Gateway Congestion Control Survey. *RFC 1254*, August 1991.
- [15] S. McCanne and S. Floyd. <http://www.nrg.ee.lbl.gov/ns/>. ns-LBNL Network Simulator, 1996.
- [16] Netperf. The Public Netperf Homepage: <http://www.netperf.org/>. The Public Netperf Homepage, 1998.
- [17] T. Ott, J. Kemperman, and M. Mathis. The Stationary Behavior of Ideal TCP Congestion Avoidance. <ftp://ftp.bellcore.com/pub/tjo/TCPwindow.ps>, August 1996.
- [18] V. Paxson. End-to-End Internet Packet Dynamics. In *Proc. of ACM SIGCOMM*, September 1997.
- [19] K. K. Ramakrishnan and S. Floyd. A Proposal to Add Explicit Congestion Notification (ECN) to IPv6 and to TCP. *Internet Draft draft-kksif-ecn-03.txt*, October 1998.
- [20] K. K. Ramakrishnan and R. Jain. A Binary Feedback Scheme for Congestion Avoidance in Computer Networks. *ACM Transaction on Computer Systems*, 8(2):158–181, May 1990.
- [21] A. Romanow and S. Floyd. Dynamics of TCP Traffic over ATM Networks. *Proc. of ACM SIGCOMM*, pages 79–88, August 1994.
- [22] B. Suter, T. V. Lakshman, D. Stiliadis, and A. Choudhury. Design Considerations for Supporting TCP with Per-flow Queueing. *Proc. IEEE INFOCOM*, March 1998.
- [23] C. Villamizar and C. Song. High Performance TCP in ANSNET. *Computer Communication Review*, 24(5):45–60, October 1994.