

# Traffic Sensitive Active Queue Management for Improved Multimedia Streaming

Vishal Phirke, Mark Claypool, and Robert Kinicki

Worcester Polytechnic Institute, Worcester, MA 01609, USA  
{claypool|rek}@cs.wpi.edu

**Abstract.** The Internet, which has traditionally supported throughput-sensitive applications such as email and file transfer, is increasingly supporting delay-sensitive multimedia applications such as interactive audio. These delay-sensitive applications would often rather sacrifice some throughput for lower delay. Unfortunately, the current Internet does not offer choices in the amount of delay or throughput an application receives, but instead provides monolithic best-effort service to all applications. This paper proposes and evaluates a new Active Queue Management (AQM) technique that employs source hints to provide service at network routers that is sensitive to the Quality of Service (QoS) expectations for a variety of applications. Applications indicate their delay or throughput sensitivity via a delay hint in their outgoing packets. The router, which we call RED-Boston, uses the delay hints to dynamically adjust the router to yield better delay performance for delay-sensitive applications and better throughput for throughput-sensitive applications. Using a new QoS metric, our simulations demonstrate that RED-Boston yields higher QoS than an adaptive version of RED for both throughput-sensitive flows and delay-sensitive flows. RED-Boston operates equally well in all traffic scenarios and fits the current best-effort Internet environment without requiring traffic monitoring.

## 1 Introduction

The Internet today carries traffic for applications with a wide range of delay and loss requirements. Traditional applications such as FTP and E-mail are primarily concerned with throughput, while Web traffic is moderately sensitive to delay as well as throughput. Emerging applications such as IP telephony, video conferencing and networked games have different requirements in terms of throughput and delay than these traditional applications. In particular, interactive multimedia applications, unlike traditional applications, have more stringent delay constraints than loss constraints. Moreover, with the use of repair techniques [3, 17, 20, 21] packet losses can be partially or fully concealed, enabling multimedia applications to operate over a wide range of losses, and leaving end-to-end delays as the major impediment to acceptable quality.

Unfortunately, current Internet routers are not able to provide applications a choice in Quality of Service (QoS).<sup>1</sup> Due to the simplicity of the FIFO queuing mechanism, drop-tail buffers are the most widely used queuing scheme in Internet routers today. When drop-tail buffers overflow, newly arriving packets are dropped regardless of the application type of the arriving packet. To accommodate bursty traffic, drop-tail routers on the Internet backbone are over-provisioned with large FIFO buffers. When faced with persistent congestion, these drop-tail routers yield high delays for all flows passing through the bottlenecked router. This best-effort service provides no consideration for multimedia flows that can be severely affected by high delays.

Typical AQM mechanisms [8, 9, 6, 2, 15, 12] provide equal treatment to incoming traffic and tend to be tuned for high throughput without specific consideration for the traversing application's delay requirements. Other AQM mechanisms [25, 16, 4, 1, 18, 7] monitor the bandwidth use by unresponsive flows but ignore the willingness of interactive multimedia applications to accept reduced throughput if accompanied by reduced delay.

ABE [13] provides a queue management mechanism for low delay traffic. In ABE, delay-sensitive applications can sacrifice throughput for lower delays. However, ABE traffic classification is rigid in that applications are either delay-sensitive or throughput-sensitive without the applications themselves being able to choose relative degrees of sensitivity to throughput and delay.

DiffServ [11, 14] and other class-based approaches [22, 5] try to give differentiated service to traffic aggregates. However, they require complicated mechanisms to negotiate service level agreements. In addition, DiffServ architectures require traffic monitors, markers, traffic shapers, classifiers and droppers to enable components to work together.

We present a new AQM approach, called *RED-Boston*, for IP networks that provides QoS service for applications with a variety of delay and throughput requirements, such as file transfer, Web browsing and interactive audio. Unlike approaches that provide fixed classes of service, each application sending traffic into the RED-Boston router chooses a customized delay-throughput trade-off based on its own requirements. The service is still best-effort in that it requires no additional policing mechanisms, charging mechanisms or usage control. Under the proposed service, RED-Boston routers operate equally well under scenarios with only traditional traffic and operate better under scenarios with mixed or multimedia only traffic.

In our approach, applications mark each packet with a suggested delay, referred to as a *delay hint*, indicating the relative importance of delay versus throughput. At a congested router, the queue manager, which we have named *RED-Boston*<sup>2</sup>, uses the delay hints of all packets to calculate an average delay

---

<sup>1</sup> Throughout this work, we use the term QoS to refer explicitly to *delay* and *throughput* provided by the network.

<sup>2</sup> Similar to various versions of TCP, which are named after cities in Nevada, we have named our version of the RED [8] active queue management technique after a city in Massachusetts.

hint and a target average queue size chosen to best fit the aggregate traffic. RED-Boston inserts packets in the outbound queue based on their delay hint relative to the average delay hint. Packets with a delay hint lower than the average delay hint are inserted towards the front of the queue, while being dropped with a higher than average drop probability. Conversely, packets with a delay hint higher than the average delay hint are inserted towards the back of the queue, while being dropped with a lower than average drop probability. Thus, using a low delay hint, delay-sensitive applications experience reduced end-to-end delay, while receiving the lower throughput resulting from an increased drop rate. Similarly, using a high delay hint, throughput-sensitive applications receive an increased throughput resulting from a decreased drop rate, while having an increased end-to-end delay.

We evaluate RED-Boston via simulation under a wide variety of traffic mixes. The results show that RED-Boston achieves better overall application performance for both delay-sensitive multimedia applications and throughput-sensitive applications. As traffic mix changes from mostly throughput-sensitive to mostly delay-sensitive, RED-Boston gives consistently better performance than Adaptive RED [9].

This paper is organized as follows: Section 2 explains the RED-Boston mechanism and provides an example of how an application may use the RED-Boston service; Section 3 describes simulation experiments used to evaluate RED-Boston; Section 4 analyzes the simulation results; and Section 5 contains conclusions and suggests possible future work.

## 2 RED-Boston

In RED-Boston, applications mark each outgoing packet with a delay hint that suggests the relative importance of delay versus throughput for the packet. Delay hints are explained in Section 2.1. The router mechanism itself is explained in Section 2.2. The router uses the delay hints to calculate a target average queue size as explained in Section 2.2.1. Packets experience a drop probability based on their delay hint and average drop probability, as explained in Section 2.2.2, while being inserted in the outgoing queue based on their delay hint relative to the average delay hint, as explained in Section 2.2.3. Lastly, the space and processing requirements for RED-Boston are discussed in Section ??.

### 2.1 Delay Hints

In RED-Boston, application end-hosts indicate their sensitivity to queuing delay as source hints in the IP packet header that we refer to as *delay hints*. The delay hint is not an absolute bound on queuing delay, but is rather a suggestion to Internet routers as to the relative importance of delay versus throughput.

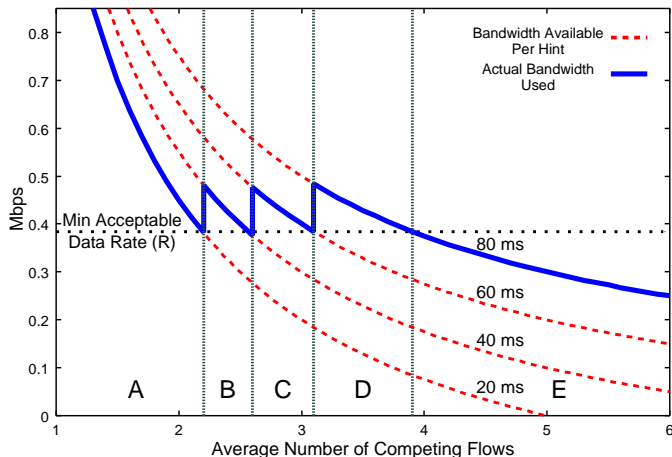
RED-Boston does not guarantee delay based on the delay hint, but the delay hint facilitates the router in providing an appropriate service. At the RED-Boston router, a relatively low delay hint indicates an application's desire for

lower delay and lower throughput while a higher delay hint suggests the application prefers higher throughput even if this implies higher delays. Applications such as FTP or Email that want to minimize overall transfer time without significant concern over individual packet delays would choose a high delay hint. Applications such as multimedia streaming that seek to minimize end-to-end delay could choose to send a low delay hint. However, under congested network conditions these applications may discover that the low delay hint also yields unacceptably low throughput. Applications could then raise their delay hints until measured throughput reached acceptable levels. Thus, sources are able to vary their delay hints in any manner that increases their user-defined QoS. Packets that carry no delay hints are handled using a default delay hint that corresponds to the router's target queue size.

Figure 1 presents a simple illustration of how a delay-sensitive TCP-friendly application might use delay hints. An interactive videoconference running over a company T1 link competes with  $n$  TCP flows for the bandwidth. The videoconference requires a minimum data rate,  $R$ , of about 384 Kbps<sup>3</sup> to insure acceptable video quality. This rate is depicted by the horizontal dotted line in Figure 1. Under these conditions, the videoconference is free to choose the source hint. The curved dashed lines in Figure 1 depict the approximate bandwidth an application would receive for source hints of 20 ms, 40 ms, 60 ms and 80 ms (reading left to right). The higher the source hint, the more bandwidth the videoconference will receive, but the higher the delay. The lower the source hint, the less bandwidth the videoconference will receive, but the lower the delay. The "best" source hint for the application depends upon the perceived quality,  $pq$ , for each throughput and delay combination. For this simplified example, we assume that  $pq$  is 0 when the data rate,  $R'$ , is less than  $R$  and increases inversely with the delay for  $R' > R$ . In other words, once the videoconference obtains its minimum data rate, the greatest benefit to perceived quality comes from lower delay. Under these conditions, the vertical lines in Figure 1 indicate regions  $A$ ,  $B$ ,  $C$ ,  $D$  and  $E$  where the videoconference would choose different delay hints. In region  $A$ , the required minimum bandwidth can be obtained by using the lowest delay hint, 20 ms. As the T1 link becomes more congested, the videoconference would choose source hints of 40 ms and 60 ms for regions  $B$  and  $C$  respectively. In region  $D$ , the videoconference is forced to use the largest source hint, 80 ms, to obtain the minimum bandwidth. This is the same source hint as is likely used by the competing TCP flows. In region  $E$ , even with a large source hint, the available bandwidth drops below the minimum acceptable rate, in which case the videoconference user would probably terminate the connection. Note that this example is greatly simplified and we would expect perceived quality functions that trade off throughput and latency in a more sophisticated manner would be used.

---

<sup>3</sup> A typical minimum data rate for an H.261 videoconference.



**Fig. 1.** Possible Delay Hint Strategy for an Interactive Multimedia Application

## 2.2 Mechanism

Arriving delay hints affect three aspects of the RED-Boston router mechanism. First, the target average queue size moves to adjust the average queue size as described in Section 2.2.1. Second, the delay hint determines the drop probability, where higher delay hints mean lower drop probabilities and vice versa, as explained in Section 2.2.2. Third, the delay hint also determines where the incoming packet is placed in the outgoing queue as described in Section 2.2.3.

**2.2.1 Moving Target** Adaptive RED [9] extends RED [8] by keeping the average queuing delay  $q_{ave}$  within a target range. In ARED, the targeted average queue size is independent of traffic QoS requirements. We argue that the targeted average queue size should be based on the incoming traffic's requirements. RED-Boston extends ARED by providing a moving target queue size based on an exponentially weighted moving average of the delay hints of the incoming packets, as shown below:

$$target = (1 - w_t)target + w_t C \times delay\_hint / (packet\_size \times 8) \quad (1)$$

where  $C$  is the router's outgoing link capacity and  $w_t$  is the weight for calculating the exponentially weighted moving average (set to  $w_t = 0.02$  in our implementation). The weighted average allows the queue target to stabilize even if the individual packets arriving from different flows have different delay hints. When overall traffic changes from throughput-sensitive to delay-sensitive, the target reflects the average queuing delay requirements of the incoming packets.

**2.2.2 Delay Hint Based Drop Probability** Since RED-Boston provides different queuing delays for different delay hints, the drop probability must be adjusted based on delay hint. RED-Boston adjusts the ARED per-packet drop probability  $p'$  based on the delay hint and the average delay as follows:

$$p' = p \times q_{delay} / delay\_hint \quad (2)$$

where  $q_{delay}$  is the queuing delay corresponding to the average queue size. The moving target described in Section 2.2.1 allows RED-Boston to maintain the average queue size at approximately the average of the aggregate traffic's delay requirements, providing fairness in the per-packet drop probability calculation.

**2.2.3 Weighted Insert** The RED-Boston outgoing queue is sorted based on packet weights. An incoming packet is inserted in the sorted queue based on its weight which is computed as shown below:

$$weight = arrival\ time + delay\_hint \quad (3)$$

Arrival time is used as an aging mechanism in this calculation to avoid packet starvation. Under normal conditions a flow's delay hints will not vary and differences in arrival times will prevent packet reordering. Since packets are served from the front of the sorted queue, the packet with the earliest deadline will be served first. Note the delay hint is not an upper bound on queuing delay for the packet but rather a relative hint such that packets with lower delay hints experience lower delays than packets with higher delay hints. However, our results show that on average, packets receive a queuing delay close to the delay hint specified.

Thus, RED-Boston's *weighted insert* provides delay-sensitive packets with a relatively lower queuing delay. On the other hand, RED-Boston's *delay hint based drop probability* compensates by providing delay tolerant packets with a lower drop probability and hence higher throughput.

**2.2.4 Algorithm** Figure 2 summarizes the RED-Boston algorithm. Each packet contains a delay hint, as described in Section 2.1. For each incoming packet, the delay hint is used to update the moving target for average queue size, kept by the router, as described in Section 2.2.1. If there is extreme congestion, indicated by the queue average being above  $max_{th}$ , the packet is dropped. If there is congestion, as indicated by the queue average being between  $min_{th}$  and  $max_{th}$ , the drop probability for the packet is computed based on the queue parameters, the *delay\_hint*, and the average queuing delay, *delay*, as described in Section 2.2.2. If the packet is not dropped, it is inserted in the queue based on a weighting of its arrival time and delay hint, as described in Section 2.2.3.

### 3 Experiments

This section discusses the experimental methods and NS-2 simulation details associated with our comparison study of the performance of RED-Boston ver-

```

on receiving packet pkt:
// Calculate moving target, (see Section 2.2.1)
target = (1 - wt) target + wt C × pkt.delay_hint / (pkt.size × 8)

if (qavg ≥ maxth) then

    dropPacket(pkt, 1)

elseif (qavg ≥ minth) then

    // Calc drop prob p, based on RED (see [8])
    p = calcDropP(qavg, minth, maxth, maxp)

    // Calc delay hint based drop prob, p' (see Section 2.2.2)
    p' = p × (qdelay / delay_hint)

    if (!dropPacket(pkt, p')) then

        // Insert in the queue based on weight, (see Section 2.2.3)
        weight = arrival time + pkt.delay
        insertPacket(pkt, weight)

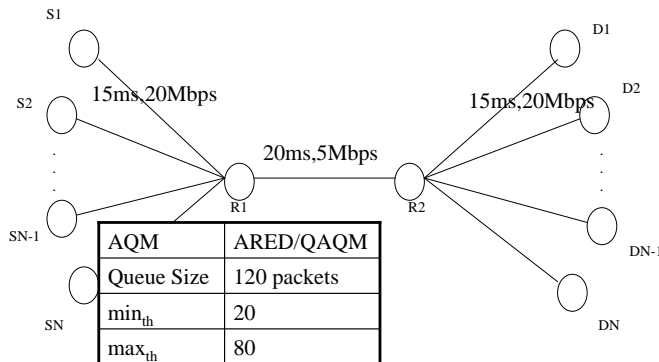
Every 500 msec: // Default in [9]
    // Adjust maxp so qave hits target, (see [9] or [23])
    maxp = adaptMaxP(maxp, target, qavg)

```

**Fig. 2.** RED-Boston Algorithm

sus ARED. The NS-2 simulator [19] provides the ability to simulate drop-tail, RED and ARED routers. Additionally, NS-2 includes code to simulate both TCP protocols and TCP-friendly rate controlled protocols such as TFRC [10]. RED-Boston was simulated by extending the ARED implementation with the *moving target*, *weighted insert* and *delay hint based drop probability* algorithms as explained in Section 2 and, TCP NewReno and TFRC were modified to send delay hints.

The generic network topology used for all experiments is shown in Figure 3. S1-SN are traffic sources and D1-DN are destinations. The S-D pairs were varied to provide traffic mixes that included different mixes of TCP NewReno and TFRC flows (see Section 4), where the TFRC flows are meant to represent TCP-friendly multimedia flows. All sources send fixed-length 1000-byte packets. All links connecting sources to router R1 and all links connecting destinations to router R2 have a 20 Mbps capacity and a 15 ms delay. With the bandwidth and delay of the bottleneck link going from R1 to R2 set at 5 Mbps and 20 ms respectively, this topology is such that packet drops only occur at R1, where all measurements are made. For both RED-Boston and ARED, the queue size at the



**Fig. 3.** Network Topology

congested router is 120 packets and  $min_{th}$  and  $max_{th}$  are set to 20 packets and 80 packets respectively. All simulated flows start at time 0 sec and end at time 100 sec. Performance measurements are taken during the stable interval between 20 sec and 80 sec to avoid transient conditions from startup and stopping.

To ascertain overall router performance, we measure link utilization, average queue size and drop rate. Moreover, since application performance realistically involves tradeoffs between throughput and delay, we believe analyzing either alone is not sufficient. We propose a new metric,  $QoS$ , that provides a quantitative performance metric designed to capture to some degree the nature of the throughput-delay tradeoff. The correct choice of QoS function depends upon the the application's sensitivity to delay and throughput. As illustrated in Section 2.1, an application may dynamically adjust its delay hints to current network conditions in attempt to improve its QoS. To include a wide spectrum of individual application types that can select their own tailored QoS objective, we provide a QoS measure that is generic while permitting customized combinations of delay and throughput measures for each individual application:

$$QoS = T^\alpha / D^\beta \quad (4)$$

where  $\alpha + \beta = 1$ . While the choices of  $\alpha$  and  $\beta$  will depend upon the throughput and delay tolerances of the specific applications themselves, the expectation is that the relativity of the delay hints from RED-Boston enabled sources will vary accordingly. Note our new metric, QoS, is in fact a traffic-sensitive variant of the power performance metric.

In the experiments discussed in the next section,  $\alpha = 1$ , is used for throughput-sensitive flows (such as Email or file transfer), since these flows can tolerate delay, and QoS for them depends only upon throughput. We use  $\alpha = 0.5$  and  $\beta = 0.5$  for delay-sensitive flows (such as an interactive videoconference), since these flows require low delay, but also moderate amounts of throughput. For medium

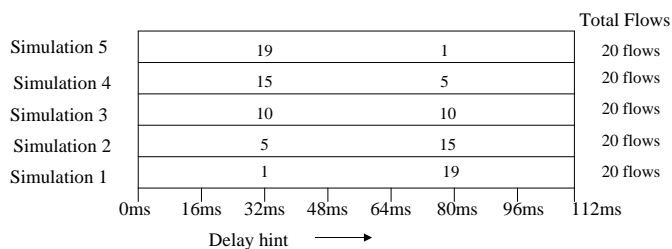


throughput and medium delay-sensitive flows (such as some HTTP flows for Web browsing), we would use  $0.5 < \alpha < 1$  and  $0 < \beta < 0.5$ .

## 4 Analysis

This section presents results analysis of one set of NS-2 simulations designed to compare RED-Boston’s performance against ARED. We focus on changing the percentages of delay-sensitive and throughput-sensitive flows in the incoming traffic mix. Results from a second set of simulations consisting of incoming traffic with a range of throughput and delay requirements is presented in [23], but is not presented here due to lack of space.

The objective of the simulations is to evaluate RED-Boston’s robustness as the incoming traffic mix at the router varies from mostly throughput-sensitive traffic to mostly delay-sensitive traffic. Figure 4 provides details on the traffic mixes used for the first five simulations.



**Fig. 4.** Simulation Traffic Mix

Each simulation ran 20 flows using the network topology shown in Figure 3. Throughput-sensitive flows were represented by TCP flows carrying a delay hint of 80 ms, which corresponds to a queue size of 50 packets at the bottlenecked router R1. Delay-sensitive flows were represented by TFRC flows carrying a delay hint of 32 ms which corresponds to a queue size of 20 packets at R1. The X-axis for all the graphs in this section indicate the changing traffic mix corresponding to the five different simulations described in Figure 4.

Figure 5 shows the average queue size in packets for RED-Boston and ARED as the traffic mix varies. While the average queue size remains constant across all the ARED simulations, the moving target mechanism in RED-Boston allows the average queue size to adjust to the delay hint distributions caused by changes in the incoming traffic’s QoS requirements. Thus, when traffic is mostly throughput-sensitive, the average queue size is higher to increase the overall throughput. Correspondingly, when traffic is mostly delay-sensitive, the average queue size becomes smaller to reduce the overall queuing delays. The RED-Boston average queue sizes in these five experiments reflects the average

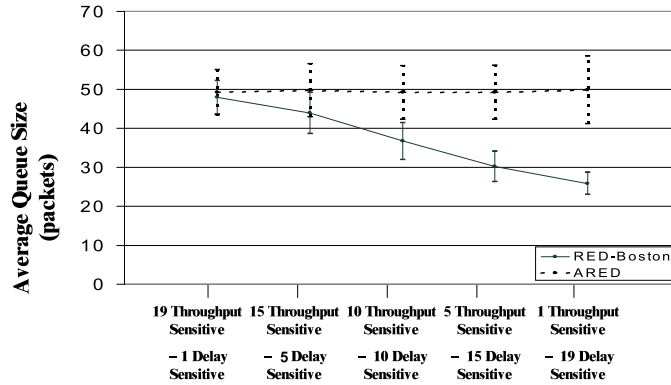


Fig. 5. Average Queue Size

requirements of incoming traffic. On the other hand, Adaptive RED's target range is insensitive to the incoming traffic's requirements.

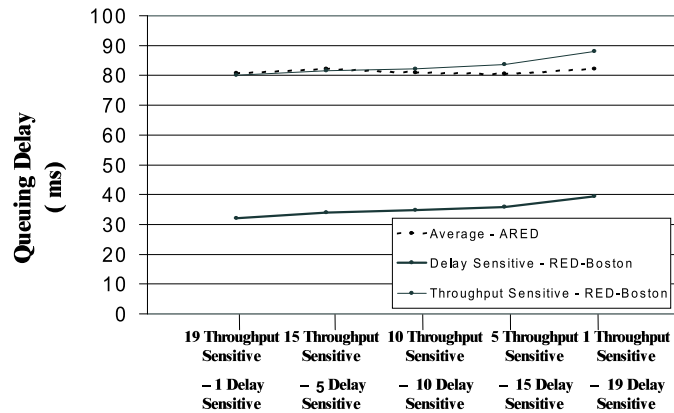


Fig. 6. Queuing Delays

Figure 6 shows the queuing delays experienced by the throughput-sensitive and the delay-sensitive flows. For ARED, graphing only the average delay curve is necessary because the throughput-sensitive and delay-sensitive flows receive identical treatment from ARED despite having different QoS requirements. Figure 6 indicates that RED-Boston gives lower queuing delays to delay-sensitive flows as compared to throughput-sensitive flows. Moreover, when the delay-sensitive flows represent a small percentage of the incoming traffic, the queuing delays experienced by the delay-sensitive flows is close to the 32 ms delay

hint. However, as the percentage of delay-sensitive flows increases in the incoming traffic mix, the average queuing delay for these flows increases slightly. This is caused by many delay-sensitive flows competing with each other for lower delays. Conversely, when most of the incoming flows are throughput-sensitive, their average queuing delay is near the 80 ms delay hint. As the percentage of delay-sensitive flows increases in the incoming traffic, the average queuing delay for throughput-sensitive flows increases slightly. This increase happens because RED-Boston lets the delay-sensitive flows cut in the queue ahead of throughput-sensitive flows to get lower delays. However, from Figure 6, it is clear that the aging component of the RED-Boston packet weight calculation prevents starvation of throughput-sensitive flows.

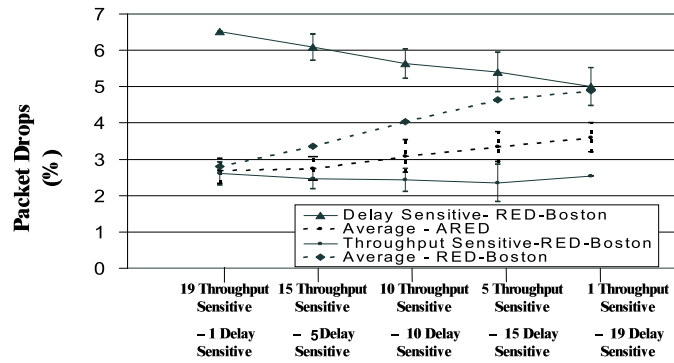


Fig. 7. Average Per Flow Percentage of Packets Dropped

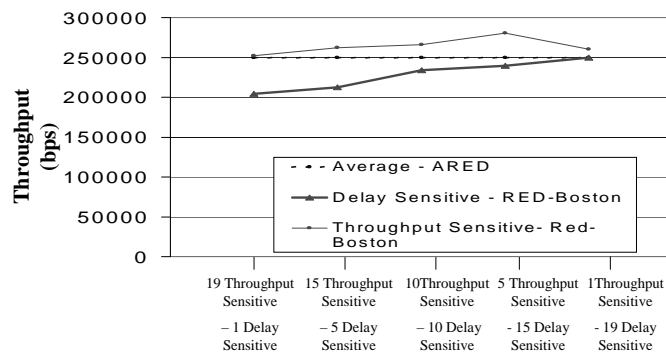


Fig. 8. Average Per Flow Throughput

Figure 7 shows the average percentage of packets dropped per flow in RED-Boston and ARED. Again only one curve is needed for ARED, since ARED treats all flows equally regardless of QoS requirements. For RED-Boston, Figure 7 does separate out the average per flow percentage of packets dropped for throughput-sensitive flows and delay-sensitive flows. The delay sensitive flow competing with 19 throughput-sensitive flows has to pay a high drop rate for the low delay service it receives because it is inserted far below the average queue size determined by the throughput sensitive flows. As the fraction of delay-sensitive flows increases, the average per flow drop rates for the delay-sensitive flows decreases. This is because as the number of delay sensitive flows increases, RED-Boston's moving target brings the average queue size down to give a lower delay service and the corresponding drop rate is shared equally by all the delay-sensitive flows. For throughput-sensitive flows, the RED-Boston drop rate is consistently lower than that of ARED and is almost constant for all simulations irrespective of number of delay-sensitive flows.

Figure 8 shows the average per flow throughput for ARED and average per flow throughput for throughput-sensitive flows and delay-sensitive flows for RED-Boston. The ARED throughput is nearly constant for all simulations and is the same for all flows. In RED-Boston, delay-sensitive flows get lower throughput than throughput-sensitive flows in exchange for low delay service. When there are 19 throughput-sensitive and one delay-sensitive flow, the delay-sensitive flow sacrifices some throughput for lower delay. This additional throughput sacrificed by the one delay-sensitive flow is equally shared by 19 throughput-sensitive flows. However, as the number of delay-sensitive flows increases, their throughput penalty decreases since the average queue size is reduced. The throughput sacrificed by delay-sensitive flows is shared equally by the existing throughput-sensitive flows in each case. In practice, as described in Section 2.1, if the reduced throughput given to a delay-sensitive flow falls below an acceptable threshold for that flow, the flow can send a larger delay hint if it is able to trade higher delay for higher throughput.

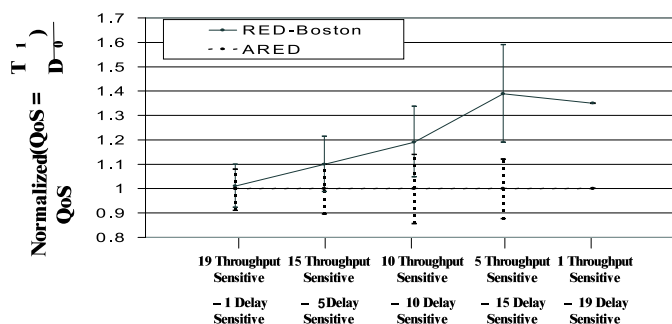
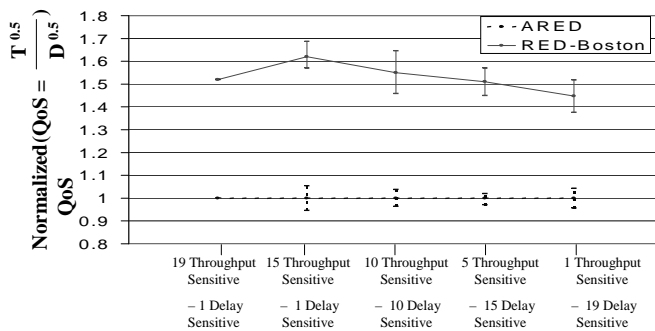


Fig. 9. Normalized QoS for Throughput-Sensitive Flows

Figure 9 presents the QoS for throughput-sensitive flows with a RED-Boston router normalized with respect to the QoS for throughput-sensitive flows traversing an ARED router. As described in Section 3, QoS for throughput-sensitive flows is defined as  $QoS = T^1/D^0$ . The RED-Boston throughput-sensitive flows record consistently higher QoS than the equivalent ARED flows. Furthermore, the RED-Boston QoS advantage grows as the percentage of delay-sensitive flows in the traffic mix increases. When delay-sensitive flows dominate the mix, more flows are willing to give up throughput for lower delay and as a result the throughput sensitive flows experience higher QoS even though they are in the minority.



**Fig. 10.** Normalized QoS for Delay Sensitive Flows

Paralleling Figure 9, Figure 10 graphs the QoS for delay-sensitive flows in RED-Boston normalized with respect to the QoS for delay-sensitive flows in ARED. As described in Section 3, the specific QoS used for delay-sensitive flows is  $QoS = T^{0.5}/D^{0.5}$ . Using this metric, RED-Boston yields a QoS improvement (on average) for delay-sensitive flows more than 50 percent over delay-sensitive flows served by an Adaptive RED router. The QoS for delay-sensitive flows in RED-Boston decreases as the number of delay-sensitive flows increases, because they start competing with each other for low delay service. The resulting improvement seen in the QoS for the delay sensitive flows captures the improvement in the quality that interactive multimedia applications using TFRC would experience with RED-Boston routers.

## 5 Summary

This paper presents RED-Boston, an active queue management technique designed to provide better QoS support to all applications under the current best-effort Internet environment. With RED-Boston, applications include a delay hint in their packets as an indication of their preference for either lower queuing delay

or higher throughput at the router. RED-Boston gives a relatively lower delay service to delay-sensitive flows and relatively higher throughput to throughput-sensitive flows. Thus, not only can low-bandwidth multimedia flows, such as interactive audio, obtain the low delays they need to support interactive communication, but delay-tolerant non-interactive flows, such as file downloading, can obtain the bandwidth they require to support efficient network transfer.

RED-Boston does not provide the delay and throughput guarantees provided by IntServ approaches, but nor does RED-Boston have the inherent scalability difficulties required by such per-flow reservations. RED-Boston does not pre-define classes of service as do DiffServ approaches, but instead provides a continuum of service classes. The service afforded by RED-Boston does not require additional monitoring or charging of either flows that take advantage of the low delay service or flows that take advantage of the high throughput service.

In evaluating RED-Boston with varying percentages of delay-sensitive and throughput-sensitive flows, we find that RED-Boston is able to adapt the overall router queue parameters to better suit the current traffic mix. Additionally, RED-Boston provides differentiated services by providing lower delay for the delay-sensitive flows and higher-throughput for throughput-sensitive flows. This improves QoS support for all flows in RED-Boston as compared to ARED.

Currently, RED-Boston expects flows to be responsive to network congestion in a TCP-friendly fashion. Since RED-Boston provides a higher drop rate for flows with a low delay hint, unresponsive flows should not gain more of a bandwidth advantage under RED-Boston than they do under current Internet environments. Still, an extension to RED-Boston would be to enhance it with a rate-based active queue management technique such as CSFQ [25], so that the dropping probability can be calculated based on a delay hint and the flow rate, helping ensure bandwidth fairness for unresponsive traffic.

Another extension to RED-Boston could be interaction between cascaded RED-Boston routers. When dequeuing a packet, RED-Boston could be modified to update the cumulative amount of time the packet has waited in router queues, thus providing additional information for trading of delay and throughput for RED-Boston routers downstream.

## References

1. A. Clerget and W. Dabbous. Tag-based Unified Fairness. In *Proceedings of IEEE INFOCOM*, April 2001.
2. S. Athuraliya, V. H. Li, S. H. Low, and Q. Yin. REM: Active Queue Management. *IEEE Network*, May 2001.
3. J-C. Bolot, S. Fosse-Parisis, and D. Towsley. Adaptive FEC-Based Error Control for Internet Telephony. In *Proceedings of IEEE INFOCOM*, March 1999.
4. Z. Cao, Z. Wang, and E. Zegura. Rainbow Fair Queuing: Fair Bandwidth Sharing Without Per-Flow State. In *Proceedings of IEEE INFOCOMM*, March 2000.
5. Jae Chung and Mark Claypool. Dynamic-CBT and ChIPS - Router Support for Improved Multimedia Performance on the Internet. In *Proceedings of the ACM Multimedia Conference*, November 2000.

6. W. Feng, D. Kandlur, D. Saha, and K. Shin. Blue: An Alternative Approach To Active Queue Management. In *Proceedings of the Workshop on Network and Operating Systems Support for Digital Audio and Video (NOSSDAV)*, June 2001.
7. W. Feng, D. Kandlur, D. Saha, and K. Shin. Stochastic Fair Blue: A Queue Management Algorithm for Enforcing Fairness. In *Proceedings of IEEE INFOCOM*, April 2001.
8. S. Floyd and V. Jacobson. Random Early Detection Gateways for Congestion Avoidance. *IEEE/ACM Transactions on Networking*, August 1993.
9. Sally Floyd, Ramakrishna Gummadi, and Scott Shenker. Adaptive RED: An Algorithm for Increasing the Robustness of RED's Active Queue Management. Under submission, 2001. <http://www.icir.org/floyd/papers/adaptiveRed.pdf>.
10. Sally Floyd, Mark Handley, Jitendra Padhye, and Jorg Widmer. Equation-Based Congestion Control for Unicast Applications. In *Proceedings of ACM SIGCOMM Conference*, pages 45 – 58, 2000.
11. J. Heinanen, F. Baker, W. Weiss, and J. Wroclawski. Assured Forwarding PHB Group. *IETF Request for Comments (RFC) 2597*, June 1999.
12. C. Hollot, V. Misra, D. Towsley, and W. Gong. On Designing Improved Controllers for AQM Routers Supporting TCP Flows. In *Proceedings of IEEE INFOCOM*, April 2001.
13. P. Hurley, M. Kara, J. Le Boudec, and P. Thiran. ABE: Providing a Low Delay within Best Effort. *IEEE Network Magazine*, May/June 2001.
14. V. Jacobson, K. Nichols, and K. Poduri. Expedited Forwarding PHB Group. *IETF Request for Comments (RFC) 2598*, June 1999.
15. S. Kunniyur and R. Srikant. Analysis and Design of an Adaptive Virtual Queue. In *Proceedings of ACM SIGCOMM*, August 2001.
16. D. Lin and R. Morris. Dynamics of Random Early Detection. In *Proceedings of ACM SIGCOMM Conference*, September 1997.
17. Yanlin Liu and Mark Claypool. Using Redundancy to Repair Video Damaged by Network Data Loss. In *Proceedings of IS&T/SPIE/ACM Multimedia Computing and Networking (MMCN)*, January 2000.
18. Debasis Mitra, Keith Stanley, Rong Pan, Balaji Prabhakar, and Konstantinos Psounis. CHOKe, A Stateless Active Queue Management Scheme for Approximating Fair Bandwidth Allocation. In *Proceedings of IEEE INFOCOM*, March 2000.
19. University of California Berkeley. The Network Simulator - ns-2. Internet site <http://www.isi.edu/nsnam/ns/>.
20. C. Padhye, K. Christensen, and W. Moreno. A New Adaptive FEC Loss Control Algorithm for Voice Over IP Applications. In *Proceedings of IEEE International Performance, Computing and Communication Conference*, February 2000.
21. K. Park and W. Wang. QoS-Sensitive Transport of Real-Time MPEG Video Using Adaptive Forward Error Correction. In *Proceedings of IEEE Multimedia Systems*, pages 426 – 432, June 1999.
22. Mark Parris, Kevin Jeffay, and F. Smith. Lightweight Active Router-Queue Management for Multimedia Networking. In *Proceedings of Multimedia Computing and Networking (MMCN), SPIE Proceedings Series*, January 1999.
23. Vishal Phirke, Mark Claypool, and Robert Kinicki. Traffic Sensitive Active Queue Management for Improved Multimedia Streaming. Technical Report WPI-CS-TR-02-10, Worcester Polytechnic Institute, April 2002.
24. William Pugh. Skip Lists: A Probabilistic Alternative to Balanced Trees. *Communications of the ACM*, 33(6):668–676, June 1990.

25. Ion Stoica, Scott Shenker, and Hui Zhang. *Core-Stateless Fair Queueing: Achieving Approximately Fair Bandwidth Allocations in High Speed Networks*. In *Proceedings of ACM SIGCOMM Conference*, September 1998.
26. Ion Stoica and Hui Zhang. *Providing Guaranteed Services Without Per Flow Management*. In *Proceedings of ACM SIGCOMM Conference*, September 1999.