

Monitoring and controlling QoS network domains

By Ahsan Habib[†], Sonia Fahmy^{*†} and Bharat Bhargava

Increased performance, fairness, and security remain important goals for service providers. In this work, we design an integrated distributed monitoring, traffic conditioning, and flow control system for higher performance and security of network domains. Edge routers monitor (using tomography techniques) a network domain to detect quality of service (QoS) violations—possibly caused by underprovisioning—as well as bandwidth theft attacks. To bound the monitoring overhead, a router only verifies service level agreement (SLA) parameters such as delay, loss, and throughput when anomalies are detected. The marking component of the edge router uses TCP flow characteristics to protect ‘fragile’ flows. Edge routers may also regulate unresponsive flows, and may propagate congestion information to upstream domains. Simulation results indicate that this design increases application-level throughput of data applications such as large FTP transfers; achieves low packet delays and response times for Telnet and WWW traffic; and detects bandwidth theft attacks and service violations. Copyright © 2004 John Wiley & Sons, Ltd.

Ahsan Habib is a Postdoctoral Researcher in the School of Information Management and Systems (SIMS), University of California at Berkeley. He works with Prof. John Chuang in the areas of Peer-to-Peer Networking, Network Economics, Multihoming, and Next Generation Network Architecture. His research interests also include Network Security, Network Monitoring, Quality of Service, and Distributed Systems. He received his Ph.D. from the Department of Computer Sciences at Purdue University in August 2003. He received a M.S. from the Department of Computer Science at Virginia Tech in 1999, and a B.Sc. from Department of Computer Science and Engineering at Bangladesh University of Engineering and Technology (BUET) in 1996.

Sonia Fahmy received her PhD degree at the Ohio State University in August 1999. Since then, she has been an assistant professor at the Computer Sciences department at Purdue University. Her research interests are in the design and evaluation of network architectures and protocols. She is currently investigating Internet tomography, overlay networks, network security, and wireless sensor networks. Her work is published in over 50 papers, including publications in IEEE/ACM Transactions on Networking, Computer Networks, IEEE INFOCOM, IEEE ICNP, and ACM NOSSDAV. She received the National Science Foundation CAREER award in 2003, the Schlumberger foundation technical merit award in 2000 and 2001, and the OSU presidential fellowship for dissertation research in 1998. She is a member of the ACM, IEEE, Phi Kappa Phi, Sigma Xi, and Upsilon Pi Epsilon. Please see <http://www.cs.purdue.edu/homes/fahmy/> for more information.

Bharat Bhargava's research involves both theoretical and experimental studies in distributed systems. He has conducted experiments in large-scale distributed systems, communications, authentication, key management, fault-tolerance and Quality of Service. He was the chairman of the IEEE Symposium on Reliable and Distributed Systems held at Purdue in October 1998. He is on the editorial board of three international journals. In the 1988 IEEE Data Engineering Conference, he and John Riedl received the best paper award for their work on 'A Model for Adaptable Systems for Transaction Processing.' He is a fellow of Institute of Electrical and Electronics Engineers and Institute of Electronics and Telecommunication Engineers. He has been awarded the charter Gold Core Member distinction by IEEE Computer Society for his distinguished service. He received Outstanding Instructor Awards, from the Purdue chapter of the ACM in 1996 and 1998. He has received IEEE Technical Achievement award for a major impact of his decade long contributions to foundations of adaptability in communication and distributed systems in 1999. His students have received best paper awards in International conferences and have started a Nasdaq listed company.

*Correspondence to: Sonia Fahmy, Department of Computer Sciences, Purdue University, West Lafayette, IN 47907-2066, USA

[†]E-mail: fahmy@cs.purdue.edu

1. Introduction

The success of the Internet has increased its vulnerability to misuse and performance problems. Internet service providers are now faced with the challenging task of continuous monitoring of their network to ensure that security is maintained, and customers are obtaining their agreed-upon service. Service providers must also identify when their networks need to be re-configured or re-provisioned, and must obtain the best performance they can out of these networks. Based upon these requirements, our primary objective in this work is twofold: (i) achieving higher user-perceivable quality of service (QoS) and overall resource utilization in a network domain, and (ii) flagging under-provisioning problems and network misuse. Towards this objective, we design edge routers that combine (a) low-overhead monitoring, together with (b) traffic conditioning at network domain edges, and (c) unresponsive flow control, in order to mitigate misuse, congestion, and unfairness problems in Internet domains. Monitoring network activity has the additional benefit of detecting denial of service (DoS) and bandwidth theft attacks, which have become an expensive problem in today's Internet.

The success of the Internet has increased its vulnerability to misuse and performance problems.

Our integrated monitoring, conditioning and control techniques will be illustrated on the differentiated services (diff-serv) architecture.¹ Diff-serv is a simple approach to enhance quality of service (QoS) for data and multimedia applications in the Internet. In diff-serv, complexity is pushed to the boundary routers of a network domain to keep core routers simple. The edge routers at the boundary of an administrative domain shape, mark, and drop traffic if necessary. The operations are based on Service Level Agreements (SLAs) between adjacent domains. The SLA defines what levels of service a user can expect from a service provider. The SLA clarifies the goals of both parties, and ensures that both of them abide by the agreement.

Our design comprises three primary components. The monitoring component infers possible attacks and SLA violations. The traffic conditioning component marks TCP traffic, using basic knowledge of TCP operation. The unresponsive flow control component regulates traffic entering a domain, and conveys congestion information to upstream domains. These three edge router components, and the flow of data and control among them, are depicted in Figure 1. Our focus when designing each component will be on scalability and low overhead. We now outline the operation of each component.

—SLA Monitoring Component—

QoS-enabled networks (e.g., differentiated services networks) are vulnerable to different types of attacks from traditional IP network domains. For example, users may inject or re-mark traffic with high QoS requirements, which may cause other users to have lower throughput, or higher delay and packet loss. We define an *attack* to be an incident when a user violates his/her SLA or re-marks packets to steal bandwidth. We need to flag such SLA violations or bandwidth theft attacks. Towards this end, we will extend and exploit network tomography techniques that use end-to-end measurements to infer internal domain behavior. Although network tomography has witnessed a flurry of research activity in the past few years, these new tomography results have not been integrated with the more mature research on monitoring and control. In contrast, we will use network tomography techniques such as per-segment loss inference mechanisms² to monitor and control network domains.

—Traffic Conditioning Component—

At the edge of a network domain, traffic conditioners can utilize knowledge of TCP characteristics to give priority to 'critical' packets, and mitigate TCP bias to flows with short round trip times (RTTs). Such intelligent conditioning functions, however, have traditionally required maintaining per-flow state information. While edge routers between a stub domain and a transit domain do not typically handle very large numbers of flows, many edge routers, such as Internet Exchange points among peering domains,

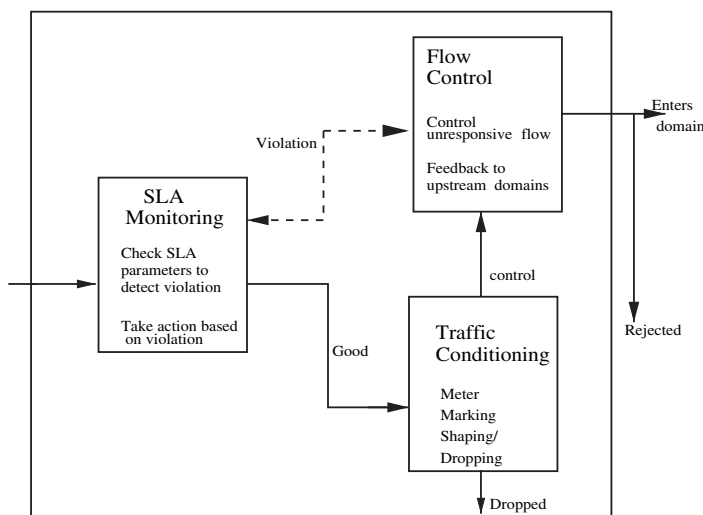


Figure 1. Monitoring, conditioning, and flow control components in an edge router

are highly loaded. To address this problem, we will design a conditioner that only uses packet header information instead of stored state when possible, and employs replacement policies to control the amount of state maintained.

—Congestion Control Component—

Congestion collapse from undelivered packets is an important problem in the Internet.³ Congestion collapse occurs when upstream bandwidth is consumed by packets that are eventually dropped downstream. This can be caused by unresponsive flows that do not reduce their transmission rates in response to congestion. The congestion collapse problem can be mitigated using improved packet scheduling or active queue management,^{4,5} but such open loop techniques do not affect congestion caused by unresponsive flows in upstream domains. To address this important problem, we will design a mechanism to control the rate at which packets enter a network domain to the rate at which packets leave the domain. Congestion is detected when many high priority packets are being dropped.⁶ Edge routers which detect or infer such drop can therefore regulate unresponsive flows. This results in better performance for users, and better overall resource utilization in the network domain.

We conduct a series of simulation experiments to study the behavior of all three components of

this framework. Our simulation results show that TCP-aware edge router marking improves throughput of greedy applications like large FTP transfers, and achieves low packet delays and response times for Telnet and WWW traffic. We also demonstrate how attacks and unresponsive flows alter network delay and loss characteristics, and hence can be detected by our monitoring component, and controlled by the congestion control component.

The remainder of this paper is organized as follows. Section 2 gives an overview of the differentiated services architecture—which we use as an underlying quality of service (QoS) framework—and discusses previous results related to the components of our proposed edge routers. Our network monitoring and loss inference techniques for attack detection are discussed in Section 3. Section 4 discusses the design of the TCP-aware traffic conditioning component. Section 5 explains how to detect and control unresponsive flows during congestion. Our simulation set-up for performance evaluation is described in Section 6. Section 7 discusses our simulation results, and conclusions are given in Section 8.

2. Background and Related Work

As previously mentioned, we use the diff-serv framework as the underlying QoS approach. In

diff-serv networks, traffic enters a domain at an *ingress* router and leaves a domain at an *egress* router. An ingress router is responsible for ensuring that the traffic entering the domain conforms to the SLA with the upstream domain. An egress router may perform traffic conditioning functions on traffic forwarded to a peering domain. In the core of the network, Per Hop Behaviors (PHBs) achieve service differentiation. The current diff-serv specification defines two PHB types: Expedited Forwarding⁷ and Assured Forwarding (AF).⁸ AF provides four classes (queues) of delivery with three levels of drop precedence (DP0, DP1, and DP2) per class. The Differentiated Services Code Point (DSCP), contained in the IP header DSFIELD/ToS field, is set to mark the drop precedence. When congestion occurs, packets marked with higher precedence (e.g., DP2) must be dropped first.

Providing QoS in diff-serv networks has been extensively studied in the literature. Clark and Fang introduced RIO in 1998,⁹ and developed the Time Sliding Window (TSW) tagger. They show that sources with different target rates can achieve their targets using RIO even for different Round Trip Times (RTTs), whereas simple RED routers cannot. Assured Forwarding is studied by Ibanez and Nichols in Reference 10. They use a token bucket marker and show that target rates and TCP/UDP interaction are key factors in determining throughput of flows. Seddigh, Nandy and Pieda¹¹ also show that the distribution of excess bandwidth in an over-provisioned network is sensitive to UDP/TCP interactions. Lin, Zheng and Hou¹² propose an enhanced TSW profiler, but their solution requires state information to be maintained at core routers.

In the next three subsections, we discuss work related to network monitoring and tomography, traffic conditioning, and congestion collapse solutions, which comprise the three components of our proposed design.

—2.1. Network Tomography and Violation Detection—

Since bottleneck bandwidth inference techniques such as packet pairs were proposed in the early 1990s, there has been increased interest in inference of internal network characteristics (e.g.,

per-segment delay, loss, bandwidth, and jitter) using correlations among end-to-end measurements. This problem is called *network tomography*. Recently, Duffield *et al.*² have used unicast packet 'stripes' (back-to-back probe packets) to infer link-level loss by computing packet loss correlation for a stripe at different destinations. This work is an extension of loss inference with multicast traffic, e.g., References 13, 14. We develop a tomography-based, low-overhead method to infer delay, loss, and throughput and detect problems that alter the internal characteristics of a network domain. The main contribution of the monitoring component of our work is to utilize existing network tomography techniques in novel ways. We also extend these tomography techniques for QoS networks.

In addition to tomography work, a number of network monitoring techniques have been recently proposed in the literature. In efficient reactive monitoring,¹⁵ global polling is combined with local event driven reporting to monitor IP networks. Breitbart *et al.*¹⁶ use probing-based techniques where path latencies and bandwidth are measured by transmitting probes from a single point of control. They find the optimal number of probes using vertex cover solutions. Recent work on SLA validation¹⁷ uses a histogram aggregation algorithm to detect violations. The algorithm measures network characteristics like loss ratio and delay on a hop-by-hop basis and uses them to compute end-to-end measurements. These are then used in validating the end-to-end SLA requirements. All of these techniques involve core routers in their monitoring. In contrast, our work pushes monitoring responsibilities to the edge routers. We use an Exponential Weighted Moving Average (EWMA) for delay, and an average of several samples for loss as in RON,¹⁸ since it is both flexible and accurate.

—2.2. Traffic Conditioning—

Edge routers perform traffic conditioning and control functions. The edge router may alter the temporal characteristics of a stream to bring it into compliance with a traffic profile specified by the network administrator.¹ A traffic meter measures and sorts packets into precedence levels. Marking, shaping, or dropping decisions are based upon the measurement result.

Marking—Markers can mark packets deterministically or probabilistically. A probabilistic packet marker, such as Time Sliding Window marker,¹⁹ obtains the current flow rate, *measuredRate*, of a user from the meter. The marker tags each packet based on the *targetRate* from the SLA and the current flow rate. An incoming packet is marked as *IN* profile (low probability to drop) if the corresponding flow has not reached the target rate, otherwise the packet is marked as high drop precedence with probability $1 - p$, where p is given by, equation (1):

$$p = \frac{\text{measuredRate} - \text{targetRate}}{\text{measuredRate}} \quad (1)$$

Shaping/dropping—Shaping reduces traffic variation and provides an upper bound for the rate at which the flow traffic is admitted into the network. A shaper usually has a finite-size buffer. Packets may be discarded if there is insufficient space to hold the delayed packets. Droppers drop some or all of the packets in a traffic stream in order to bring the stream into compliance with the traffic profile. This process is known as *policing* the stream.

A number of recent studies explore the design of more sophisticated traffic conditioners. Fang *et al.*¹⁹ proposed the Time Sliding Window Three Color Marker (TSW3CM), which we use as a standard traffic conditioner. Adaptive packet marking²⁰ uses a Packet Marking Engine (PME), which can be a passive observer under normal conditions, but becomes an active marker at the time of congestion. Yeom and Reddy²¹ also convey marking information to the sender, so that it can slow down its sending rate in the case of congestion. This requires modifying the host TCP implementation, which is difficult to deploy. Feroz *et al.*²² propose a TCP-Friendly marker. The marker protects small-window flows from packet loss by marking their traffic as *IN* profile. In this paper, we develop similar intelligent conditioning techniques *with lower overhead*.

An important problem with TCP is its bias to connections with short round trip times. Nandy *et al.* design RTT-aware traffic conditioners²³ which adjust packet marking based on RTTs, to mitigate TCP RTT bias. Their conditioner is based on the steady state TCP behavior as reported by Mathis

*et al.*²⁴ Their model, however, does not consider time-outs which we consider in this paper.

—2.3. Congestion Collapse—

As previously discussed, congestion collapse occurs when upstream bandwidth is consumed by packets that are eventually dropped downstream. A number of solutions have been proposed to mitigate this problem in Internet domains. Seddigh *et al.*²⁵ propose separating TCP (responsive to congestion) and UDP (may be unresponsive) to control congestion collapse caused by UDP. Albuquerque *et al.*²⁶ propose a mechanism, Network Border Patrol, where border routers monitor all flows, measure ingress and egress rates, and exchange per-flow information with all edge routers periodically. The scheme is elegant, but its overhead is high. Chow *et al.*²⁷ propose a similar framework, where edge routers periodically obtain information from core routers, and adjust conditioner parameters accordingly. Aggregate-based Congestion Control (ACC) detects and controls high bandwidth aggregate flows.²⁸

In the Direct Congestion Control Scheme (DCCS),⁶ packet drop of packets with the *lowest drop priority* is tracked by core routers. Dropping packets with lowest drop priority is an indication that there is severe congestion in the network. The core routers send the load information to the edge routers only during congestion. We employ a similar methodology for detecting congestion and controlling unresponsive flows. However, our proposed approach has lower storage overhead and includes a mechanism to avoid congestion collapse. Core router participation is optional.

In the next three sections, we discuss the three components of our proposed edge router design.

3. Tomography-based Violation Detection Component

An attacker can impersonate a legitimate customer of a service provider by spoofing its identity. Network filtering²⁹ can detect spoofing if the attacker and the impersonated customer are in different domains, but the attacks may proceed unnoticed otherwise. QoS domains support

low-priority classes, such as best effort, which are not controlled by edge routers. The service provider should ensure that high-priority customers are getting their agreed-upon service, so that the network can be re-configured or re-provisioned if needed, and attackers which bypass or fool edge controls are prevented. In case of distributed DoS attacks, flows from various ingress points are aggregated as they approach their victim. Monitoring can control such high bandwidth aggregates at the edges, and propagate attack information to upstream domains^{30,*}.

We employ network tomography—an approach to infer the internal behavior of a network purely based on end-to-end measurements. We use an edge-to-edge measurement-based loss inference technique to detect service violations and attacks in a QoS domain. The measurements do not involve any core router in order to scale well. We measure SLA parameters such as delay, packet loss, and throughput to ensure that users are obtaining their agreed upon service. Delay is defined as the edge-to-edge latency; packet loss is the ratio of total flow packets dropped in the domain[†] to the total packets of the same flow which entered the domain; and throughput is the total bandwidth consumed by a flow inside a domain. If a network domain is properly provisioned and no user is misbehaving, the flows traversing the domain should not experience excessive delay or loss. Although jitter (delay variation) is another important SLA parameter, it is flow-specific and therefore, not suitable to use in network monitoring. In this section, we describe edge-to-edge inference of delay, loss and throughput, and a violation detection mechanism.

—3.1. Delay Measurements—

Delay bound guarantees made by a provider network to customer flows are for the delays experienced

by the flows between the ingress and egress edges of the provider domain. For each packet traversing an ingress router, the ingress copies the packet IP header into a new packet with a certain pre-configured probability p_{probe} . The ingress encodes the current time into the payload and marks the protocol identifier field of the IP header with a new value. The egress router recognizes such packets and removes them from the network. Additionally, the egress router computes the packet delay for flow i by subtracting the ingress time from the egress time. (We assume NTP is used to synchronize the clocks.) The egress then sends the packet details and the measured delay to an entity we call the *SLA monitor* which typically resides at an edge router. At the monitor, the packets are classified as belonging to customer j , and the average packet delay of the customer traffic is updated using an exponential weighted moving average (EWMA) (we use a current sample weight 0.2). If this average packet delay exceeds the delay guarantee in the SLA, we conclude that this may be an indication of an SLA violation.

—3.2. Loss Inference—

Packet loss guarantees made by a provider network to a customer are for the packet losses experienced by its conforming traffic inside the provider domain. Measuring loss by observing packet drop at all core routers and communicating them to the SLA monitor at the edge imposes significant overhead. We use packet stripes² to infer link-level loss characteristics inside the domain. A series of probe packets with no delay between the transmission of successive packets, or what is known as a ‘stripe’, is periodically transmitted. For a two-leaf tree spanned by nodes $0, k, R_1, R_2$, stripes are sent from the root 0 to the leaves to estimate the characteristics of three links (Figure 2). For example, the first two packets of a 3-packet stripe are sent to R_2 and the last one to R_1 . If a packet reaches a receiver, we can deduce that the packet has reached the branch point k . By monitoring the packet arrivals at R_1, R_2 and both, we can write equations with three known quantities and estimate the three unknown quantities (loss rates of links $0 \rightarrow k, k \rightarrow R_1$ and $k \rightarrow R_2$) by applying conditional probability definitions, as discussed in Reference 2. We combine estimates of several

**As with any detection mechanism, the attackers can attack the mechanism itself, but the cost to attack our distributed monitoring mechanism is higher than the cost to inject or spoof traffic, or bypass a single edge router.*

†A flow can be a micro-flow identified by source and destination addresses and ports and protocol identifier, or an aggregate of several micro-flows.

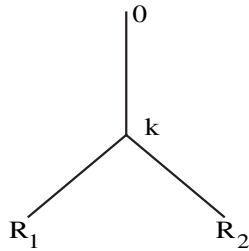


Figure 2. Binary tree to infer per-segment loss from source 0 to receivers R_1 and R_2

stripes to limit the effect of non-perfect correlation among the packets in a stripe. This inference technique extends to trees with more than 2 leaves and more than 2 levels.²

We extend this end-to-end unicast probing scheme to routers with multiple active queue management instances, e.g., 3-color RED,³¹ and develop heuristics for the probing frequency and the particular receivers to probe to ensure good domain coverage. Assured forwarding queues use three drop precedences referred to as green, yellow, and red. Suppose P_{red} is the percentage of red packets accepted (not dropped) by the active queue. We define percentages for yellow and green traffic similarly, and show how these percentages are computed in the Appendix. Link loss can be inferred by subtracting the transmission probability from 1. If L_g , L_y , and L_r are the inferred losses of green, yellow and red traffic, respectively, loss can be expressed as:

$$L_{\text{class}} = \frac{n_g P_{\text{green}} L_g + n_y P_{\text{yellow}} L_y + n_r P_{\text{red}} L_r}{n_g + n_y + n_r} \quad (2)$$

where n_i is the number of samples taken from i -colored packets. However, when loss of green traffic is zero, we take the average of yellow and red losses. When the loss of yellow traffic is also zero, we report only loss of red probes. We reduce the overhead of loss inference by probing the domain links with high delay only, as determined by the delay measurement procedure. We also measure throughput by probing egress routers *only* when delay and loss are excessive. This helps pinpoint the user or aggregate which is consuming excessive bandwidth, and causing other flows to receive lower quality than their SLAs.

—3.3. Violation and Attack Detection—

When delay, loss, and bandwidth consumption exceed the pre-defined thresholds, the monitor concludes there may be an SLA violation or attack. Excessive delay is an indication of abnormal conditions inside the network domain. If there are losses for the premium traffic class, or if the loss ratios of assured forwarding traffic classes exceed certain levels, a possible SLA violation is flagged. The violation can be caused by aggressive or unresponsive flows, denial of service attacks, flash crowds, or network under-provisioning. To detect distributed DoS attacks, the set of links with high loss are identified. If high bandwidth aggregates traversing these high-loss links have the same destination IP prefix, there is either a DoS attack or a flash crowd, as discussed in Reference 28.

When delay, loss, and bandwidth consumption exceed the pre-defined thresholds, the monitor concludes there may be an SLA violation or attack.

Decisions are taken by consulting the destination entity. Jung *et al.* analyze the characteristics of flash crowds and DoS attacks in Reference 32. Their work reveals several distinguishable features between these two. For example, the client distribution of a flash crowd event follows popular distribution among ISPs and networks, however, this is not true for a DoS attack. The other distinguishable features are per client request rate, overlap of clusters a site sees before and after the event, and popularity distribution of the file accessed by the clients. Using these characteristics, the monitor can decide whether it is a DoS attack or a flash crowd. If this is determined to be an attack, the appropriate ingress routers are notified and the offending user traffic is throttled, as discussed in Section 5.

4. Traffic Conditioning Component

We incorporate several techniques in the traffic conditioning component to improve performance

of applications running on top of TCP. The key idea behind these techniques is to protect critical TCP packets from drop, without requiring large state tables to be maintained. We protect SYN packets (as indicated in the TCP header) by giving them low drop priority (DP0).[‡] Since TCP grows the congestion window exponentially until it reaches the slow start threshold, *ssthresh*, and the congestion window is reduced to 1 or half of the *ssthresh* for time-outs or packet loss, we also protect small window flows from packet losses by marking them with DP0, as proposed in Reference 22. Edge routers use sequence number information in packet headers in both directions to determine if the window is small. ECN-Capable TCP may reduce its congestion window due to a time-out, triple duplicate ACKs, or in response to explicit congestion notification (ECN).³³ In this case, TCP sets the CWR flag in the TCP header of the first data packet sent after the window reduction. Therefore, we give low drop priority to a packet if the CWR or ECN bit is set. This avoids consecutive *ssthresh* reductions that lead to poor performance with TCP Reno.³⁴ We also mark packets inversely proportionally to the square of the flow requested rates if proportional sharing of excess bandwidth is required.²³ The marker avoids marking high drop priority in bursts, in order to work well with TCP Reno, as proposed in Reference 22.

An optional feature to mitigate the TCP short RTT bias is to mark based on RTT and RTO, if such information is available to the edge router. To understand how RTT and RTO information can be used to increase fairness, we review two TCP throughput models. Equation (3) shows that, in a simple TCP model that considers only duplicate ACKs,³⁵ bandwidth is inversely proportional to RTT where *MSS* is the maximum segment size and *p* is the packet loss probability:

$$BW \propto \frac{MSS}{RTT \sqrt{p}} \quad (3)$$

An RTT-aware marking algorithm based on this model (proposed in Reference 23) works well for

[‡]With SYN DoS attacks, this may be unfavorable. Therefore, if a mechanism to detect SYN attacks is included in the router, this option can be turned off as soon as a SYN attack is detected.

a small number of flows because equation (3) accurately represents the fast retransmit and recovery behavior when *p* is small. We have observed that for a large number of flows, short RTT flows time out because only long RTT flows are protected by the conditioner after satisfying the target rates. To mitigate this unfairness, we use the throughput approximation by Padhye *et al.*³⁵:

$$BW \propto \frac{MSS}{RTT \sqrt{\frac{2bp}{3}} + T_o \times \min\left(1, 3\sqrt{\frac{3bp}{8}}\right) p(1 + 32p^2)} \quad (4)$$

where *b* is the number of packets acknowledged by a received ACK, and *T_o* is the time-out length. Designing an RTT-aware traffic conditioner using equation (4) is more accurate than using equation (3) because it considers time-outs. Simplifying this equation, we compute the packet drop ratio between two flows, ρ as:

$$\rho^2 = \left(\frac{RTT_1}{RTT_2}\right)^2 \times \frac{T_{o1}}{T_{o2}} \quad (5)$$

where *RTT_i* and *T_{o_i}* are the RTTs and timeouts, respectively, of flow *i*.³⁶ If the measured rate is beyond the target rate of a flow, the marker marks the packets as DP0 with probability $(1 - p^2)$ where *p* is defined in equation (1). The unmarked packets are out-of-profile (DP1 and DP2) packets. These packets are directly related to the packet drop probabilities at the core. This means that packet drop at the core is proportional to the out-of-profile marked packets. Equation (5) is used to mark out of profile packets as DP1 or DP2, where such packets will be dropped before DP0 during congestion. The RTT and RTO are estimated at the edge routers using an algorithm similar to that specified in Reference 23.

Each of the techniques discussed in this section has advantages and limitations. Protecting SYN, ECN, and CWR packets, and marking according to the target rate do not need to store per flow information and are simple to implement. On the other hand, protecting small window flows and marking according to the RTT and RTO values requires maintaining and processing per flow information. To bound state overhead at the edge routers, we store per flow information at the edge only for a certain number of flows based on available memory. The edge router uses a least recently used (LRU) state replacement policy when the


```

for For each incoming flow do
  if there is a complete state entry for this flow then
    statePresent = TRUE
    Update the state table
  else
    statePresent = FALSE
    Add the flow in the state table (replace if needed)
  end if
  if statePresent is TRUE then
    Use Standard conditioner plus SYN, ECN, CWR, small window, burst, RTT-RTO based marking
  else
    Use Standard conditioner plus SYN, ECN, and CWR based marking
  end if
end for

```

Figure 3. Algorithm for adaptive traffic conditioner

number of flows exceeds the maximum number that can be maintained. Therefore, for every flow, conditioning is based on state information if it is present. If there is no state present, conditioning only uses techniques that rely on header information. The conditioner pseudo-code is given in Figure 3.

5. Congestion Control Component

This section describes the congestion control component of the edge router. This component detects and controls unresponsive flows, i.e., flows that do not reduce their rates as a response to congestion. As previously discussed, SLA monitors can inform edge routers of congestion inside a domain. A shaping algorithm at the edge routers can therefore control unresponsive flows at the time of congestion. In addition, ingress routers of a domain may propagate congestion information to the *egress routers of upstream domains*. A stub domain that is connected to only end-users does not propagate this information to the users, instead, it controls the incom-

ing flows to conform with the service level agreements.

—5.1. Optional Core Router Detection—

In Section 3, we have shown how tomography-based loss inference techniques can be applied to detect per-segment losses using edge-to-edge probes. An alternative strategy is to track excessive drop of only the high priority (i.e., green or DP0) packets at core routers, as proposed in Reference 6. The core router maintains the tuple {source address, destination address, source port, destination port, protocol identifier, timestamp, *btlnkbw*} for dropped DP0 packets. The outgoing link bandwidth at the core, *btlnkbw*, helps regulate the flow. Edge routers shape more aggressively if the core has a 'thin' outgoing link. The core sends this drop information to the ingress routers *only* when the total drop exceeds a local threshold (thus the flow appears to be non-adaptive). The information is sent only during congestion, and only for the flows that are violating SLAs. Thus, this design does not significantly reduce the scalability of the differentiated service framework.

—5.2. Metering and Shaping—

At the egress routers, we distinguish two types of drops: a drop due to metering and shaping at downstream routers (*sdrop*), and a drop due to congestion (*cdrop*) (either obtained via inference as discussed in Section 3 or communicated from core to edge routers as in Section 5.1). For a particular flow, assume the bottleneck bandwidth is *btlnkbw* (as given above); the bandwidth of the outgoing link of the flow at the ingress router is *linkbw*; the flow has an original profile (target rate) of *targetrate*; and the current weighted average rate for this flow is *wavg*. In case of *cdrop*, the profile of the flow is updated temporarily (to yield rate *newprofile*) using equations (6) and (7) where $0 < \gamma < 1$ is the congestion control aggressiveness parameter:

$$\text{decrement} = \text{cdrop} \times \text{packet_size} \times \max\left(1, \gamma \frac{\text{linkbw}}{\text{btlnkbw}}\right) \quad (6)$$

$$\text{newprofile} = \max(0, \min(\text{newprofile} - \text{decrement}, \text{wavg} - \text{decrement})) \quad (7)$$

A higher value of γ speeds up convergence, but application QoS may deteriorate. A lower value makes traffic smoother, but it takes longer to readjust the rate. The ‘max’ term in the equation can be ignored if the bottleneck bandwidth information cannot be obtained (i.e., tools like pathchar or Net-timer cannot be used), or core router detection (Section 5.1) is unavailable. In equation (7), the weighted average of the arrival rate is computed using the Time Sliding Window [10] algorithm.

For *sdrop*, the profile is adjusted as follows:

$$\text{newprofile} = \max(0, \text{newprofile} - \text{sdrop} \times \text{packet_size}) \quad (8)$$

The *newprofile* is initialized to *targetrate*. In the absence of drops, the router increases the adjusted profile periodically at a certain rate *increment*. The rate *increment* is initialized to a constant number of packets each time the router receives drop information, and is doubled when there is no drop, until it reaches a threshold $\frac{\text{wavg}}{f}$, and then it is increased linearly. Thus, the rate adjustment algorithm follows the spirit of TCP congestion control. At the edge, shaping is based on the current

average rate and the adjusted profile. For each incoming flow, if the current average rate is greater than the adjusted profile, unresponsive flow packets are dropped.

6. Simulation Set-up

We use simulations to study the effectiveness of our edge router design. The **ns-2** simulator³⁷ with the differentiated services implementation of Nortel Networks³⁸ is used. We use the following RED parameters $\{\text{min}_{\text{thr}}, \text{max}_{\text{thr}}, P_{\text{max}}\}$: for DP0 {40, 55, 0.02}; for DP1 {25, 40, 0.05}; and for DP2 {10, 25, 0.1} (as suggested in Reference 23). w_q is 0.002 for all REDs. TCP New Reno is used with a packet size of 1024 bytes and a maximum window of 64 packets. We vary the number of micro-flows (where a micro-flow represents a single TCP/UDP connection) per aggregate from 10 to 200. We compute the following performance metrics:

Throughput.—This denotes the average bytes received by the receiver application over simulation time. A higher throughput usually means better service for the application (e.g., shorter completion time for an FTP flow). For the ISP, higher throughput implies that links are well-utilized.

Packet drop ratio.—This is the ratio of total packets dropped to the total packets sent. A user can specify for certain applications that packet drop should not exceed a certain threshold.

Packet delay.—For delay sensitive application like Telnet, the packet delay is a user metric.

Response time.—This is the time between sending a request to a Web server and receiving the response back from the server.

7. Simulation Results

The objective of this preliminary set of experiments is to evaluate the effectiveness of the three components of our edge router. In the next few sections, we study the performance of each component under various conditions.

—7.1. Detecting Attacks and SLA Violations—

In this section, we investigate the accuracy and effectiveness of the delay, loss, and throughput approximation methods for detecting violations discussed in Section 3. We use a similar network topology to that used in Reference 2 as depicted in Figure 4. The violation detection mechanism works for higher link capacities and for a larger network network topology as shown in Figure 4. Having links with capacity higher than 10Mbps requires more traffic to simulate an attack. Multiple domains are connected to all edge routers through which flows enter the network domain. The flows are created from domains attached to $E1$, $E2$, $E3$, and destined to the domains connected to edge router $E6$, so that the link $C4 \rightarrow E6$ is congested. An attack is simulated on $C4 \rightarrow E6$ by flooding this link with an excessive amount of traffic entering through different ingress routers. The purpose of this simulation is to show that the edge routers can detect service violations and attacks due to flow aggregation towards a downstream domain. Many other flows are created to ensure all links carry a significant number of flows.

We first measure delay when the network is correctly provisioned or over-provisioned (and thus experiences little delay and loss). The delay of $E1 \rightarrow E6$ is 100ms; $E1 \rightarrow E7$ delay is 100ms; and $E5 \rightarrow E4$ delay is 160ms. Attacks are simulated on router $E6$ through links $C3 \rightarrow C4$ and $C4 \rightarrow E6$. With the attack traffic, the average delay of the $E1 \rightarrow E6$ link increases from 100ms to approximately 180ms. Since all the core links have a higher capacity than other links, $C4 \rightarrow E6$ becomes the most congested link.

Figure 5 shows the cumulative distribution function (CDF) of edge-to-edge delay for the link $E1 \rightarrow E6$ under various traffic loads and in the presence of attacks. When there is no attack, the end-to-end delay is close to the link transmission delay. If the network path $E1 \rightarrow E6$ is lightly loaded, for example with a 30% load, the delay does not go significantly higher than the link transmission delay. Even when the path is 60% loaded (medium load in Figure 5), the edge-to-edge delay of link $E1 \rightarrow E6$ is increased by less than 30%. Some instantaneous values of delay increase to as high as 50% of the link transmission delay, but the average value does not fluctuate too much. In both

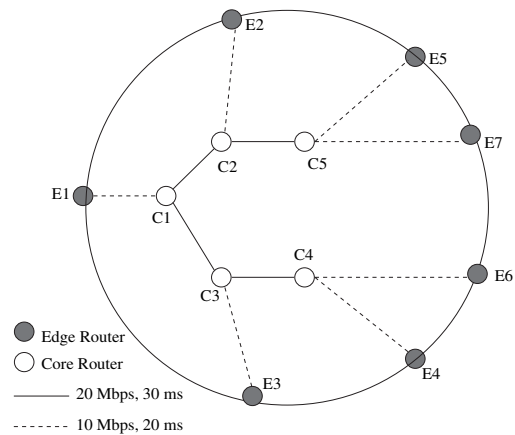


Figure 4. Topology used to infer loss and detect service violations. All edge routers are connected to multiple domains, and each domain has multiple hosts

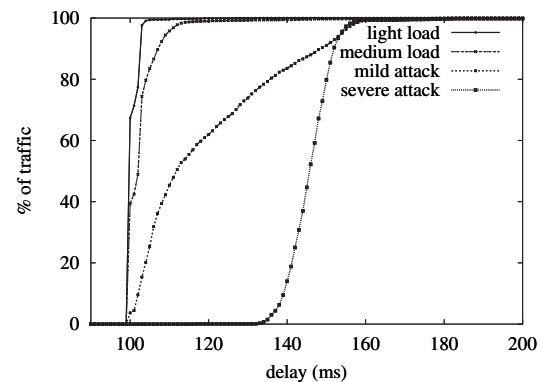


Figure 5. Cumulative distribution function (CDF) of one way delay from $E1$ to $E6$

cases, the network is properly provisioned, i.e., the flows do not violate their SLAs. On the other hand, excess traffic introduced by an attacker increases the edge-to-edge delay and most of the packets of attack traffic experience a delay 40–70% higher than the link transmission delay (Figure 5). Delay measurement is thus a good indication of the presence of excess traffic inside a network domain.

The frequency of delay probing is a critical parameter that affects the accuracy of the estimation.

Sending fewer probes reduces overhead but using only a few probes can produce inaccurate estimation, especially when some of the probes are lost in the presence of excess traffic due to an attack. We have found that sending probes at a rate of 10–15 per second is a good choice in this experiment.

We demonstrate detection of such abnormal conditions using delay measurements in three scenarios labeled 'No attack', 'Attack 1', and 'Attack 2' in Figure 6. 'No attack' indicates no significant traffic in excess of capacity. This is the baseline case of proper network provisioning and traffic conditioning at the edge routers. Attacks 1 and 2 inject more traffic into the network domain from different ingress points. The intensity of the attack is increased during time $t = 15$ seconds to $t = 45$ seconds. Loss is inferred when high delay is experienced inside the network domain. To infer loss inside a QoS network, green, yellow, and red probes are used. We use equation (2) to compute overall traffic loss per class in a QoS network. The loss measurement results are depicted in Figure 7. The loss fluctuates with time, but the attack causes packet drops of 15% to 25% in the case of Attack 1 and more than 35% with Attack 2. We find that it takes approximately 10 seconds for the inferred loss to converge to the same value as the real loss in the network. Approximately 20 stripes per second are required to infer a loss ratio close to the actual value. For more details on the probing frequencies and convergence of the estimations, see Reference 30.

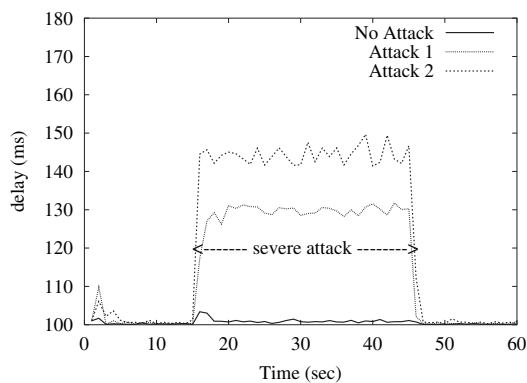


Figure 6. Observed delay on the path $E1 \rightarrow E6$ at the time of an attack. 'Attack 1' results in packet loss in excess of 15–30%. 'Attack 2' increases packet loss to more than 35%

Delay and loss estimation, together with the appropriate thresholds, can thus indicate the presence of abnormal conditions, such as distributed DoS attacks and flash crowds. When the SLA monitor detects such an anomaly, it polls edge devices for throughputs of flows. Using these outgoing rates at egress routers, the monitor computes the total bandwidth consumption by any particular user. This bandwidth consumption is then compared to the SLA bandwidth. By identifying the congested links and the egress routers connected to the congested links, the downstream domain where an attack or crowd is headed is identified. Using IP prefix matching, we determine whether many of these flows are being aggregated towards a specific network or host. If the destination confirms this is an attack, we control these flows at the ingress routers.

—7.2. Adaptive Conditioning—

As discussed in Section 4, TCP-aware marking can improve application QoS. We first perform several experiments to study each marking technique separately and study all combinations. We find that protecting SYN packets is useful for short-lived connections and very high degrees of multiplexing. Protecting connections with small window sizes (SW) contributes the most to total bandwidth gain, followed by protecting CWR packets and SYN. SW favors short RTT connections, but it reduces packet drop ratio and time-outs for long RTT connections as well, compared

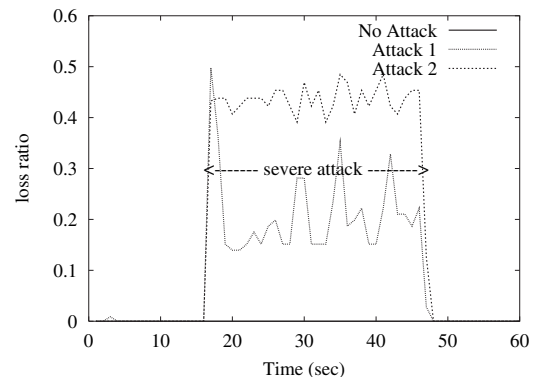


Figure 7. Packet loss ratio in link $C4 \rightarrow E6$ follows the same pattern as the edge-to-edge delay of the path $E1 \rightarrow E6$

to a standard traffic conditioner. Not marking in bursts is effective for short RTT connections. If SW is not used, Burst + CWR achieves higher bandwidth than any other combination. The RTT-RTO based conditioner mitigates the RTT-bias among short and long RTT flows. This is because when the congestion window is small, there is a higher probability of time-outs in the case of packet drops. Protecting packets (via DP0 marking) when the window is small reduces time-outs, especially back-to-back time-outs. A micro-flow also recovers from time-outs when RTO as well as RTT is used to mark packets. All these marking principles are integrated together with an adaptive state replacement policy, as given in Figure 3. We now evalu-

ate the performance of this adaptive traffic conditioner with FTP and CBR traffic, Telnet and WWW applications. The network hosts and routers are ECN-enabled for all experiments in this section, since we use the ECN and CWR packet protection mechanism. Additional results can be found in Reference 39.

Figure 9(a) compares the bandwidth with the standard and with the adaptive (Figure 3) conditioner for the simple topology shown in Figure 8(a). The total throughput is measured over the entire simulation time at the receiving end. The curve labeled 'Max' denotes the bandwidth when the standard conditioner is combined with all marking techniques and stores per-flow informa-

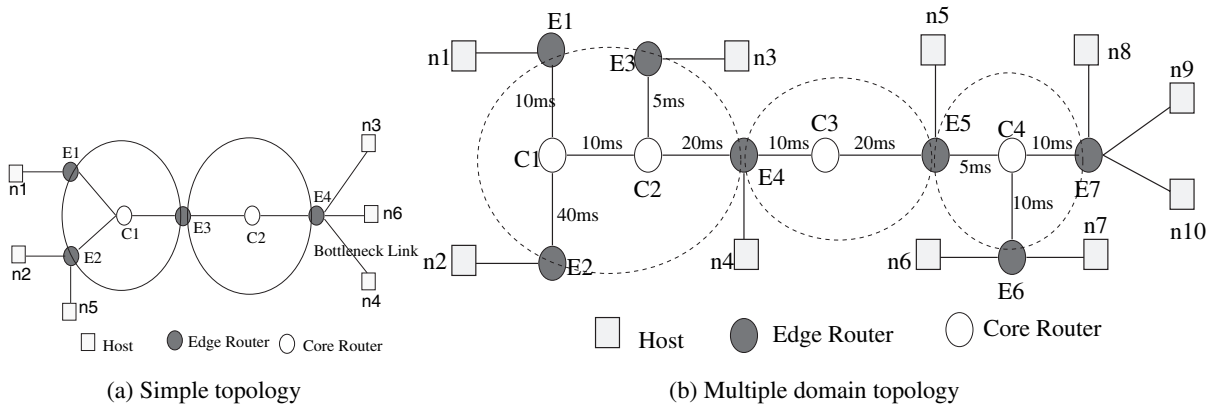


Figure 8. Simulation topologies. All links are 10Mbps. Capacity of the bottleneck links are altered in some experiments

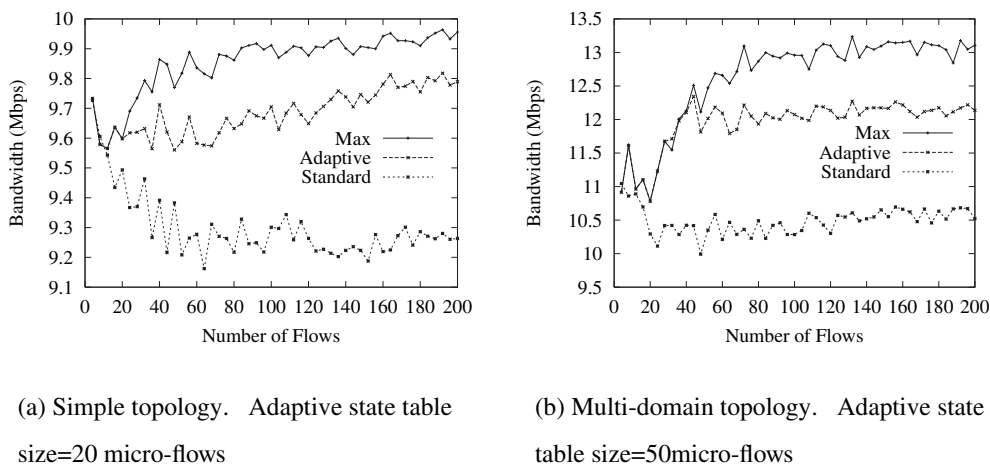


Figure 9. Achieved bandwidth with the standard conditioner and adaptive conditioner. 'Max' is the bandwidth when the standard conditioner is combined with all marking techniques and stores per-flow information for all flows

tion for all flows. We find that the adaptive conditioner outperforms the standard one for all aggregate flows. The adaptive conditioner is more fair (not shown) in the sense that short RTT flows do not steal bandwidth from long RTT flows and total achieved bandwidth is close to 10Mbps (bottleneck link speed). Therefore, with simple processing (Figure 3), performance can be increased.

Figure 8(b) depicts a more complex simulation topology where three domains are interconnected (all links are 10Mbps). The link delay between host and the edge is varied from 1 to 10ms for different hosts connected to a domain to simulate users at variable distances from the same edge routers. Aggregate flows are created between nodes $n1 \rightarrow n8$, $n2 \rightarrow n9$, $n3 \rightarrow n4$, $n5 \rightarrow n6$, and $n7 \rightarrow n9$. Thus, flows are of different RTTs, and experience bottlenecks at different links. Not all flows start/stop transmission at the same time: flows last from less than a second to a few seconds. $C2 \rightarrow E4$, $E5 \rightarrow C4$ and $C4 \rightarrow E7$ are the most congested links. Figure 9(b) shows the total bandwidth gain for this topology with different conditioners. From the figure, the adaptive conditioner performs better than the standard one, and achieves performance close to the maximum capacity. We also analyzed the bandwidth gain of each individual aggregate flow. The flows achieve similar bandwidth gains except for flows with extremely short RTTs. Thus, the adaptive conditioner improves fairness between short and long RTT flows, *without requiring large state tables*.

When each aggregate flow contains 200 micro-flows, the soft state table for the adaptive conditioner covers only a small percentage (4.16%) of the flows passing through it. We use a table for the

50 most recent micro-flows. Table 1 shows that the bandwidth achieved with the adaptive conditioner always outperforms standard conditioner. Note that when critical TCP packets are protected, they are charged from the user profile to ensure that UDP traffic is not adversely affected.

We also study performance with Telnet (delay-sensitive) and WWW (response time sensitive) applications. For the Telnet experiments, the performance metric used is the average packet delay time for each Telnet packet. We use the topology in Figure 8(b), but capacity of the $C1 \rightarrow E4$ and $E5 \rightarrow E7$ links is changed to 0.5Mbps and all other link capacities are 1Mbps to introduce congestion. We simulate 100 Telnet sessions among hosts $n1 \rightarrow n8$, $n2 \rightarrow n9$, $n3 \rightarrow n4$, $n5 \rightarrow n6$, and $n7 \rightarrow n9$. A session transfers between 10 and 35 TCP packets. Results show that the adaptive conditioner reduces packet delay over the standard conditioner for short RTT flows.

Since web traffic constitutes most (60%–80%) of the Internet traffic, we study our traffic conditioner with the WWW traffic model in **ns-2**.³⁷ Details of the model are given in Reference 40. The model uses HTTP 1.0 with TCP Reno. The servers are attached to $n6$, $n8$ and $n9$ in Figure 8(b), and $n1$, $n2$ and $n5$ are used as clients. A client can send a request to any server. Each client generates a request for 5 pages with a variable number of objects (e.g., images) per page. The default **ns-2** probability distribution parameters are used to generate inter-session time, inter-page time, objects per page, inter-object time, and object size (in kB). The network set-up is the same as with Telnet traffic. Table 2 shows the average response

Micro-flows	Standard bandwidth	Adaptive bandwidth	Max bandwidth	Adaptive (% flows covered at E4)
10	12.65	12.87	12.87	41.16
50	12.18	13.84	14.20	16.66
100	11.67	13.48	14.89	8.33
200	11.77	13.61	14.91	4.16

Table 1. Bandwidth shown is in Mbps. State table size = 50 micro-flows

Conditioner	Avg response time (sec), first pkt	Std dev	Avg response time (sec), all pkts	Std dev
Standard	0.48	0.17	2.23	0.78
Adaptive	0.45	0.14	2.15	0.75

Table 2. Response time for WWW traffic. Number of concurrent sessions = 50

time per WWW request received by the client. Two response times are shown in the table; one is to get the first packet and another is to get all data. The table shows that our adaptive conditioner reduces response time over the standard traffic conditioner. The adaptive conditioner does not change the response time significantly if the network is not congested.

The adaptive conditioner does not change the response time significantly if the network is not congested.

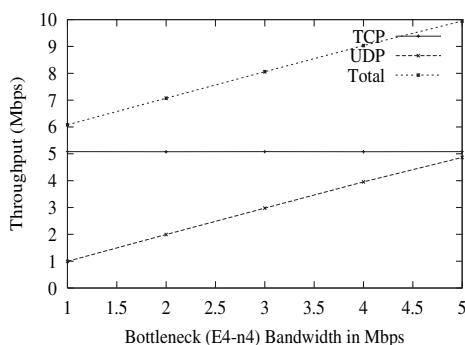
—7.3. Congestion Control—

We conduct experiments to demonstrate the role of the congestion control mechanism in preventing congestion collapse. Figure 8(a) depicts the simple topology used to demonstrate congestion collapse due to unresponsive flows. An aggregate TCP flow with 10 micro-flows from host $n1 \rightarrow n3$ and a UDP aggregate flow with 10 micro-flows from host $n2 \rightarrow n4$ are created. Both flows have the same target rate (5Mbps). Figure 10 shows how TCP and UDP flows behave with respect to changing the bottleneck bandwidth ($btlnkbw$) from 1–5Mbps. The x -axis denotes the $btlnkbw$ and the y -axis gives the throughput achieved by both flows. Figure 10(a) shows that the TCP flow gets its share of 5Mbps all the time because it does not go through the

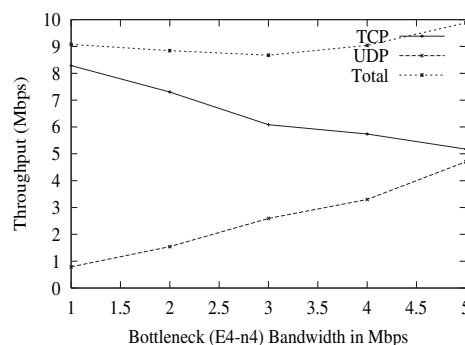
congested link. When the bottleneck bandwidth is 1Mbps, 4Mbps are wasted by UDP flows in the absence of the flow control. Figure 10(b) shows that, with flow control, the TCP flow gets an extra 8Mbps when $btlnkbw$ is 1Mbps. The flow control mechanism prevents congestion collapse due to undelivered packets.

We also experiment with varying the rate ratio, $R_r = \frac{\text{sendingRate}}{\text{Profile}}$ for UDP traffic. A value for R_r of 0.5 means that the flow is sending at 50% of its profile and a value for R_r of 4 means the flow is sending at four times its profile. When the UDP sending rate is zero, TCP can use the entire 10Mbps, and there is no shaping (shaping drop is zero) at the edge. When the UDP sending rate causes drops at the bottleneck link (e.g., when $btlnkbw = 1\text{Mbps}$), congestion collapse occurs in the absence of flow control. With flow control, even when R_r is 4 (the profile is 5Mbps and UDP is sending at 20Mbps), there is no congestion collapse.

A more complex topology with multiple domains (Figure 8(b)) and with cross traffic is also used to study the flow control framework. An aggregate of TCP flows $F1$ between $n1 \rightarrow n8$ is created, in addition to several UDP flows such as $F2$, $Cr1$, $Cr2$, and $Cr3$ between $n2 \rightarrow n9$, $n3 \rightarrow n4$, $n5 \rightarrow n6$, and $n7 \rightarrow n10$, respectively. These Cr s are used as cross traffic. The start and the finish times of the Cr flows are set differently to change the overall traffic load over the path for the flows $F1$ and $F2$. There are 10 micro-flows per aggregate in



(a) No Flow Control



(b) With Flow Control

Figure 10. (a) Without flow control and the TCP gets only 5Mbps when bottleneck bandwidth is 1Mbps of topology in Figure 8(a). (b) With flow control and now the TCP gets 8Mbps. Both flows have the same profile

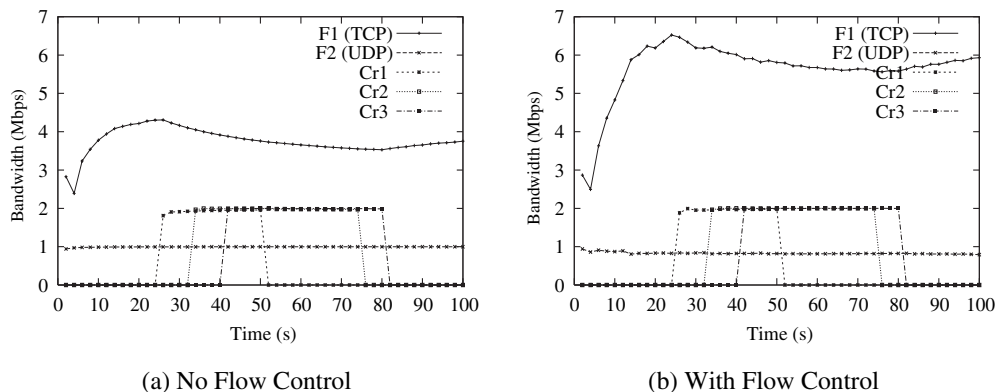


Figure 11. Dynamic adjustment of $F2$ flow works fine in the presence of cross traffic. TCP flow ($F1$) gets more bandwidth with the flow control scheme

this set-up. Flows $F1$ and $F2$ have the same profile with target rate 5Mbps, and the cross traffic sending rate is 2Mbps.

Figure 11 illustrates the bandwidth of these aggregate flows with and without flow control. The cross traffic achieves the same target in both schemes, because the flows do not encounter any bottleneck. If there is no flow control, $F1$ (TCP) cannot achieve its target of 5Mbps. With flow control, $F1$ obtains more than the target. This is because, after controlling UDP, TCP uses the remaining bandwidth.

8. Conclusions

We have investigated the design of edge routers that include tomography-based edge-to-edge probing methods to detect service level agreement violations in QoS networks, together with TCP-aware conditioning and flow control for unresponsive flows. In addition to the primary goal of this design, which is to increase performance of compliant flows, the proposed mechanisms are useful for network re-dimensioning, as well as for detecting and controlling flooding-based attacks. We have designed methods that use edge-to-edge packet stripes to infer loss for different drop precedences in a QoS network, based on observed delays. Aggregate throughputs are then measured to detect distributed denial of service attacks or flash crowds.

Marking, shaping, and policing are also adapted to respond to detection results and adapt to flow

characteristics. We give priority to critical TCP packets and mark according to flow characteristics. We use an adaptive conditioner that overwrites previous state information based on a least recently used strategy. Marking is based on information in packet headers if state information for a flow is unavailable. The adaptive conditioner is shown to improve FTP throughput, reduce packet delay for Telnet, and response time for WWW traffic. The conditioner also mitigates TCP RTT bias if it can deduce the flow RTT and RTO. Finally, we have designed a simple method to regulate unresponsive flows to prevent congestion collapse due to undelivered packets.

Most of our mechanisms can be adapted to other architectures that support service differentiation, or to active queue management techniques at network routers. For example, the RED algorithm at network routers can itself protect critical TCP packets, e.g., CWR marked packets, from drop without requiring any additional state. The adaptive conditioner concept can also be employed to keep some window size information and use that in RED dropping decisions.

Acknowledgements

The authors would like to thank Srinivas R. Avasarala and Venkatesh Prabhakar for many discussions and ideas about the monitoring component. This research is supported in part by the National Science Foundation grants ANI-0238294 (CAREER), ANI-0219110, CCR-001712 and CCR-001788, CERIAS; an IBM SUR grant, the Purdue

Research Foundation, and the Schlumberger Foundation technical merit award.

References

- Blake S, Black D, Carlson M, Davies E, Wang Z, Weiss W. An architecture for differentiated services. RFC 2475, December 1998.
- Duffield NG, Lo Presti F, Paxson V, Towsley D. Inferring link loss using striped unicast probes. In *Proc. IEEE INFOCOM*, Alaska, April 2001.
- Floyd S, Fall K. Promoting the use of end-to-end congestion control in the Internet. *IEEE/ACM Transactions on Networking* 1999; 7(6):458–472.
- Braden B, et al. Recommendations on queue management and congestion avoidance in the internet. RFC 2309, April 1998.
- Ott T, Lakshman T, Wong L. SRED: Stabilized RED. In *Proc. IEEE INFOCOM*, March 1999.
- Wu H, Long K, Cheng S, Ma J. A Direct Congestion Control Scheme for Non-responsive Flow Control in Diff-Serv IP Networks. Internet Draft, draft-wuht-diffserv-dccs-00.txt, August 2000.
- Jacobson V, Nichols K, Poduri K. An Expedited Forwarding PHB. RFC 2598, June 1999.
- Heinanan J, Baker F, Weiss W, Wroclawski J. Assured Forwarding PHB Group. RFC 2597, June 1999.
- Clark D, Fang W. Explicit allocation of best effort packet delivery service. *IEEE/ACM Transactions on Networking* 1998; 6(4):362–374.
- Ibanez J, Nichols K. Preliminary simulation evaluation of an assured service. draft-ibanez-diffserv-assured-eval-00.txt, August 1998.
- Seddigh N, Nandy B, Piedad P. Bandwidth assurance issues for TCP flows in a differentiated services network. In *Proc. Globecom 99*, 1999.
- Lin W, Zheng R, Hou J. How to make Assured Services more Assured. In *Proc. ICNP*, October 1999.
- Adams A, et al. The use of end-to-end multicast measurements for characterizing internal network behavior. *IEEE Communications Magazine* 2000; 38(5):152–159.
- Cáceres R, Duffield NG, Horowitz J, Towsley D. Multicast-based inference of network-internal loss characteristics. *IEEE Transactions on Information Theory*, November 1999; 45(7):2462–2480.
- Dilman M, Raz D. Efficient reactive monitoring. In *Proc. IEEE INFOCOM*, Alaska, April 2001.
- Breitbart Y, Chan CY, Garofalakis M, Rastogi R, Silberschatz A. Efficiently monitoring bandwidth and latency in IP networks. In *Proc. IEEE INFOCOM*, Alaska, April 2001.
- Chan MC, Lin YJ, Wang X. A scalable monitoring approach for service level agreements validation. In *Proceedings of the International Conference on Network Protocols (ICNP)*, pages 37–48, November 2000.
- Anderson D, Balakrishnan H, Kaashoek F, Morris R. Resilient overlay networks. In *Proc. ACM Symp. on Operating Systems Principles (SOSP)*, Banff, Canada, October 2001.
- Fang W, Seddigh N, Nandy B. A time sliding window three colour marker. RFC 2859, June 2000.
- Feng WC, Kandlur D, Saha D, Shin KG. Understanding and improving TCP performance over networks with minimum rate guarantees. *IEEE/ACM Transactions on Networking* 1999; 7(2):173–186.
- Yeom I, Reddy N. Realizing throughput guarantees in a differentiated services network. In *Proc. IEEE Int. Conf. on Multimedia Comp. and Systems*, June 1999.
- Feroz A, Kalyanaraman S, Rao A. A TCP-friendly traffic marker for IP differentiated services. In *Proc. IEEE/IFIP Eighth International Workshop on Quality of Service—IWQoS*, 2000.
- Nandy B, Seddigh N, Piedad P, Ethridge J. Intelligent traffic conditioners for assured forwarding based differentiated services networks. In *Proc. IFIP High Performance Networking*, Paris, June 2000.
- Mathis M, Semke J, Mahdavi J, Ott T. The macroscopic behavior of the TCP congestion avoidance algorithm. *ACM SIGCOMM Computer Communication Review* 1997; 27(3):67–82.
- Seddigh N, Nandy B, Piedad P. Study of TCP and UDP interaction for the AF PHB, 1999.
- Albuquerque C, Vickers B, Suda T. Network border patrol. In *Proc. IEEE INFOCOM*, 2000.
- Chow H, Leon-Garcia A. A feedback control extension to differentiated services. Internet Draft, draft-chow-diffserv-fbctrl-00.pdf, March 1999.
- Mahajan M, Bellovin SM, Floyd S, Ioannidis J, Paxson V, Shenker S. Controlling high bandwidth aggregates in the network. *ACM Computer Communication Review* 2002; 32(3):62–73.
- Ferguson P, Senie D. Network ingress filtering: Defeating denial of service attacks which employ IP source address spoofing agreements performance monitoring. RFC 2827, May 2000.
- Habib A, Fahmy S, Avasarala SR, Prabhakar V, Bhargava B. On detecting service violations and bandwidth theft in QoS network domains. *Computer Communications* 2003; 26(8):861–871.
- Floyd S, Jacobson V. Random early detection gateways for congestion avoidance. *IEEE/ACM Transactions on Networking*, 1993; 1(4):397–413. <ftp://ftp.ee.lbl.gov/papers/early.ps.gz>.
- Jung J, Krishnamurthy B, Rabinovich D. Flash crowds and denial of service attacks: Characterization and implications for cdns and web sites. In *Proc.*

- World Wide Web (WWW)*, Honolulu, Hawaii, May 2002.
33. Ramakrishnan K, Floyd S. A proposal to add explicit congestion notification (ECN) to IP. RFC2481, January 1999.
 34. Fahmy S. *New TCP Standards and Flavors, in Text: High Performance TCP/IP Networking*, chapter 11, pages 264–280. Prentice Hall, October 2003.
 35. Padhye J, Firoiu V, Towsley D, Kurose J. Modeling TCP throughput: a simple model and its empirical validation. In *Proc. ACM SIGCOMM '98*, 1998.
 36. Habib A, Bhargava B, Fahmy S. A round trip time and timeout aware traffic conditioner for differentiated services networks. In *Proc. IEEE International Conference on Communication (ICC)*, New York, April 2002.
 37. McCanne S, Floyd S. Network simulator ns-2. <http://www.isi.edu/nsnam/ns/>, 1997.
 38. Shallwani F, Ethridge J, Piedad P, Baines M. Diff-Serv implementation for ns. <http://www7.nortel.com:8080/CTL/#software>, 2000.
 39. Habib A, Fahmy S, Bhargava B. Design and evaluation of an adaptive traffic conditioner in differentiated services networks. In *Proc. IEEE International Conference on Computer Communication and Networks (IC3N)*, Arizona, pages 90–95, October 2001.
 40. Feldmann A, Gilbert AC, Huang P, Willinger W. Dynamics of IP traffic: A study of the role of variability and the impact of control. In *Proc. ACM SIGCOMM*, pages 301–313, 1999. ■

If you wish to order reprints for this or any other articles in the *International Journal of Network Management*, please see the Special Reprint instructions inside the front cover.

Appendix A

—Percentages Used in Multi-Priority Loss Inference—

Figure A1 depicts the drop probabilities in RED with three drop precedences. The *red* traffic has higher drop priority than *yellow* and *green* traffic. The red traffic is dropped with a probability P_{red} when the average queue size lies between two thresholds R_{min} and R_{max} . All incoming red packets are dropped when the average queue length exceeds R_{max} . P_{yellow} and P_{green} are similar. Suppose $\alpha_G(n)$ is the probability that an incoming green packet will be accepted by the queue given that n packets are in the queue. $\alpha_Y(n)$ and $\alpha_R(n)$ are defined similarly for yellow and red traffic respectively. The α values for green packets are defines as follows:

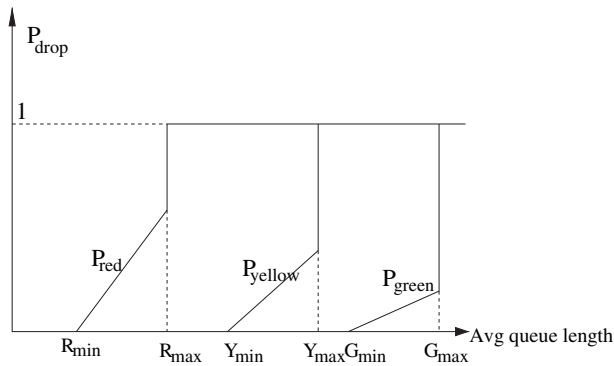


Figure A1. RED parameters for three drop precedences

$$\begin{aligned} \alpha_G(n) &= 1 && \text{if } n < G_{min} \\ \alpha_G(n) &= 0 && \text{if } n > G_{max} \\ \alpha_G(n) &= 1 - P_{green} \frac{n - G_{min}}{G_{max} - G_{min}}, && \text{otherwise} \end{aligned} \quad (A1)$$

The equations are similar for yellow and red traffic. Let P'_{red} be the percentage of packet drops due to active queue management for red packets, and let P'_{yellow} and P'_{green} be defined similarly for yellow and green, respectively. These percentages can be computed as:

$$P'_{red} = \frac{R_{max} - R_{min}}{R_{max}} \times P_{red} + \frac{G_{max} - R_{max}}{B} \times 100 \quad (A2)$$

$$P'_{yellow} = \frac{R_{max} - Y_{min}}{Y_{max}} \times P_{yellow} + \frac{G_{max} - Y_{max}}{B} \times 100 \quad (A3)$$

$$P'_{green} = \frac{G_{max} - G_{min}}{G_{max}} \times P_{green} \quad (A4)$$

where B is the buffer (queue) size at the router. The percentage of class k traffic accepted by an active queue can be expressed as:

$$P_k = 1 - P'_k \quad (A5)$$