

The War Between Mice and Elephants

LIANG GUO IBRAHIM MATTA

Computer Science Department
Boston University
Boston, MA 02215

{guo1, matta}@cs.bu.edu

May 2001

Technical Report BU-CS-2001-005

Abstract

Recent measurement based studies reveal that most of the Internet connections are short in terms of the amount of traffic they carry (mice), while a small fraction of the connections are carrying a large portion of the traffic (elephants). A careful study of the TCP protocol shows that without help from an Active Queue Management (AQM) policy, short connections tend to lose to long connections in their competition for bandwidth. This is because short connections do not gain detailed knowledge of the network state, and therefore they are doomed to be less competitive due to the conservative nature of the TCP congestion control algorithm.

Inspired by the Differentiated Services (Diffserv) architecture, we propose to give preferential treatment to short connections inside the bottleneck queue, so that short connections experience less packet drop rate than long connections. This is done by employing the RIO (RED with In and Out) queue management policy which uses different drop functions for different classes of traffic.

Our simulation results show that: (1) in a highly loaded network, preferential treatment is necessary to provide short TCP connections with better response time and fairness without hurting the performance of long TCP connections; (2) the proposed scheme still delivers packets in FIFO manner at each link, thus it maintains statistical multiplexing gain and does not misorder packets; (3) choosing a smaller default initial timeout value for TCP can help enhance the performance of short TCP flows, however not as effectively as our scheme and at the risk of congestion collapse; (4) in the worst case, our proposal works as well as a regular RED scheme, in terms of response time and goodput.

Keywords

Traffic Engineering, Congestion Control, TCP Performance, Fairness.

I. INTRODUCTION

THE ubiquitous heavy-tailed distributions in the Internet implies an interesting feature of the Internet traffic: most (e.g. 80%) of the traffic is actually carried by only a small number of connections (elephants), while the remaining large amount of connections are very small in size or lifetime (mice). In a fair network environment, short connections expect relatively fast service than long connections. However, sometimes we can not observe such nice property in the current Internet. Users may spend a long time downloading a plain text webpage — a situation for which “World-Wide Wait” was coined [1].

A careful look at the behavior of the major transport layer protocol TCP reveals some of the factors that may affect the performance of transferring a short file:

- TCP tries to conservatively ramp up its transmission rate to the maximum available bandwidth. Therefore, the sending window is initiated at the minimum possible value regardless of what is available in the network.

This work was supported in part by NSF grants CAREER ANI-0096045 and MRI EIA-9871022.

- As a reliable transport protocol, TCP couples error control with congestion control. Packet loss is detected by either the expiration of a retransmission timer or the arrival of a sequence of duplicate acknowledgments. Usually a retransmission after timeout can severely degrade the transmission rate of TCP. For a short connection, since most of the time the congestion window has a relatively small value, it does not have enough packets to activate the duplicate ACK mechanism, and thus packet loss always requires a timeout to detect.
- More importantly, TCP relies on its own packet samples to estimate an appropriate retransmission timeout (RTO) value. For the first control packets (SYN, SYN-ACK), and the first data packet, since no sampling data is available, TCP has to use a conservatively estimated initial timeout (ITO) value as RTO¹. Losing these packets can have a disastrous effect on short connection performance due to the large timeout period.
- For these reasons, short TCP flows are generally more conservative than long flows and thus tend to get less than their fair share when they compete for the bottleneck bandwidth.

In this paper, we propose to give preferential treatment to short flows² with help from an Active Queue Management (AQM) policy inside the network. We also rely on the proposed Differentiated Services (Diffserv) architecture [3] to classify flows into short and long at the edge of the network. More specifically, we maintain the length of each active flow (in packets³) at the edge routers and use it to classify incoming packets. At the core routers, we employ the RIO queue management policy [4] to isolate packets from short flows and reduce the losses experienced by these packets.

Through extensive simulations, we confirm that preferential treatment is necessary to guarantee prompt response to short TCP flows. Moreover, since our threshold-based classification method treats the initial packets from a long flow as packets from a short flow, all flows benefit from experiencing less drops of the first few packets (including control packets), which otherwise would cause unnecessary stall of the session. We also confirm that our proposed scheme can achieve better fairness and response time for short flows than schemes without preferential treatment even when the value of the ITO is reduced to match the actual configuration, as suggested in [2]. Since the RIO policy aims at maximizing statistical multiplexing, our proposed scheme does not hurt the utilization of the network. On the contrary, it sometimes can achieve better network goodput than traditional Drop Tail or RED [5] policies since it effectively avoids some unnecessary backoff of the sessions. In addition, RIO guarantees ordered delivery of packets at each link, and thus eliminates the possibility of packet reordering which may occur with other isolation-based policies (e.g. CBQ [6]).

Related Work: In a preliminary work [7], we study the interaction between short and long flows through a control theoretical approach. We then propose to isolate short and long TCP flows to enhance response time and fairness of short flows. Our results also reveal that class-based flow isolation combined with threshold-based classification at the edge may cause packet reordering, which may severely degrade TCP performance. Moreover,

¹ The recommended ITO value is 3 seconds and is implemented in most modern operating systems [2]. Some old systems use a more conservative value of 6 seconds.

² We use flow, connection and session interchangeably in this paper.

³ Without loss of generality we assume fixed-size packets and deal with packets rather than bytes.

the bandwidth (load) control in [7] is performed at each core router, while in this paper we push bandwidth control to the edges of the network. Ng *et al.* [8] explore the idea of preferential treatment of *non-TCP* bursty flows over long CBR flows to guarantee different QoS service profiles. Nandy *et al.* [9] and Firoiu *et al.* [10] study service differentiation between delay-sensitive and throughput-sensitive *long-lived* flows through active queue management.

Analytical results of short TCP flow latency under different network conditions are given in [11] and [12]. These results reveal the vulnerability of short TCP flows. Seddigh *et al.* [2] show the negative impact of the initial timeout value on the short TCP flow latency and propose to reduce the default recommended value. To enhance the performance of short TCP flows, proposals like [13], [14], [15], [16] suggest to either use a large initial window value or share the network measurement information from previous records or communication with neighbors. Notice that these proposals require modification of the TCP protocol and may lead to congestion collapse. Our proposed scheme differs from these proposals because our control is placed inside the network, and thus is more fair to clients who are unaware of such changes.

Crovella *et al.* [17, 18] explore the heavy-tailedness nature of the job size distribution and comment that size-aware job scheduling helps enhance the response time of short jobs without hurting the performance of long jobs. Their work inspires us to perform size-based service differentiation in the Internet for TCP flows.

Many studies have been conducted in the area of active queue management (AQM), most of them are variants of RED [5] or CBQ [6]. In our proposed scheme, we require the internal bottleneck routers to execute RIO (RED with In and Out) policy [4], which is a pioneering work by Clark and Feng.

Our Contribution: Our study provides an alternative usage of the Differentiated Services (Diffserv) architecture in [4, 19]. The goal in [4, 19] is to provide *long-lived* flows with “expected” bandwidth or fairness guarantees. Thus, packets are classified at the edge of the network as In or Out for that purpose, which may lead to suboptimal utilization of network resources [20]. On the other hand, our goal is to provide a *new* better-than-best-effort TCP service that attempts to *enhance the competence of short TCP flows*. This in turn creates a more fair network environment and better utilizes network resources, especially for Web traffic, the dominant traffic in the current Internet. We classify packets as In or Out to distinguish the length of the flow as Short or Long in order to preferentially treat short TCP flows inside the network.

In Section II, we show in detail the necessity and benefit of preferential treatment of short TCP flows. We then present the network architecture and mechanisms required to implement our proposed scheme in Section III. Section IV uses simulations to confirm the advantage of our scheme. In Section V, we discuss the issue of incremental deployment and extensions to the architecture. Section VI concludes this paper with future work.

II. ANALYZING SHORT TCP FLOW PERFORMANCE

In this section, we use a similar approach as in [11] to study the latency of short TCP flows. We then argue that, to achieve transmission rates comparable to that of long connections, short TCP connections should experience much less packet drops, which requires preferential treatment at the bottleneck link queues.

A. Deriving Short TCP Flow Latency

For ease of analysis, we assume packet loss happens independently from individual flow behavior. This assumption is true when the number of active flows is large. We also assume Reno-family TCP flows, however, most of the analysis can also apply to other versions of TCP. Due to space limitation, we omit a detailed description of the TCP congestion control algorithm. We refer the reader to [21] for an excellent description. Briefly speaking, TCP uses a window-based control mechanism. It employs the slow-start algorithm to quickly ramp up to the desired transmission rate, then migrates to additive increase multiplicative decrease (AIMD) control to oscillate its window around the desired value. Once severe congestion happens, it temporarily shuts down the transmission and conservatively (exponentially) estimates the time required to resume transmission. To match the *ns* simulation model, we focus our analysis on modeling the actual data packet transmission of a TCP session, which implies prompt and reliable delivery of all the control packets (SYN, SYN-ACK, etc.). Notice that such assumption favors schemes that do not give preferential treatment to such packets (see discussion in Section V). We also assume fixed-size packets and a total file size of d packets.

First Packet: As mentioned earlier, to provide reliable delivery, TCP relies on either the reception of duplicate acknowledgments (dupacks) or expiration of a retransmission timer to detect packet losses. For the first data packet, it is not possible to activate a dupack event⁴. Thus, loss of the first packet always requires at least one retransmission timeout to detect.

On the other hand, TCP uses its own packet samples to estimate the appropriate retransmission timer value. However, for the first data packet, there is no packet sample available for such purpose. Thus TCP uses a conservatively estimated value⁵ as the initial timeout (ITO) for retransmission.

Assume packet loss happens uniformly with probability p , and no loss happens on the reverse path (no ACK drops), then we can compute the expected time for the successful delivery of the first packet (assuming p is small enough so it takes at most 6 times to backoff the timer):

$$E[T_f] = RTT \cdot (1 - p) + ITO \cdot (1 - p) \sum_{i=0}^6 p^{i+1} 2^i \quad (1)$$

where the first term represents time for successful delivery without any packet loss, and the second term represents expected waiting time when loss happens.

Remaining Packets: Once the first packet is successfully delivered, TCP gains enough knowledge to estimate the retransmission timeout (RTO). The derivation in [11] can then be directly applied here (omitting setup latency) to compute the latency of transmitting the remaining data. For simplicity, we assume non-delayed acknowledgments⁶. The transmission time is given by:

$$E[T_r] = E[T_{ss}] + E[T_{loss}] + E[T_{ca}] \quad (2)$$

⁴ Unless TCP itself is modified to use a large initial window (larger than 3), however, we only model standard TCP behavior here.

⁵ 3 seconds according to [22], while some old BSD systems use 6 seconds [2].

⁶ Note that ignoring the effect of delayed ACKs on transmission latency does not qualitatively change our final conclusion.

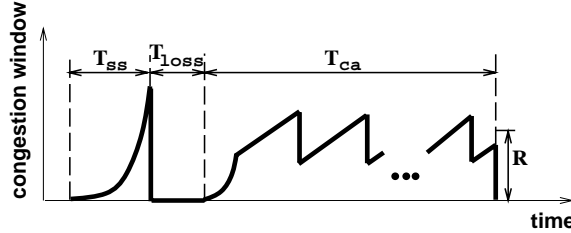


Fig. 1. TCP transmission time illustrated

where $E[T_{ss}]$ is the latency of the initial slow-start phase, $E[T_{loss}]$ deals with the packet loss that terminates the slow-start phase, and $E[T_{ca}]$ is the transmission time of the packets after TCP enters the congestion avoidance phase, as illustrated in Figure 1. These latencies are given by⁷:

$$\begin{aligned}
 E[T_{ss}] &= RTT \cdot \log_2(E[d_{ss}] + 1) \\
 E[T_{loss}] &= (1 - (1 - p)^d) \cdot (Q(p, E[W_{ss}]) \cdot E[Z^{TO}] + (1 - Q(p, E[W_{ss}])) \cdot RTT) \\
 E[T_{ca}] &= (d - E[d_{ss}]) / R(p, RTT, RTO)
 \end{aligned}$$

where $E[d_{ss}]$ is the expected number of packets transmitted in the initial slow-start phase, $E[W_{ss}]$ is the expected congestion window size reached after the initial slow-start phase, $Q(\cdot)$ is the probability of retransmission timer expiration, $E[Z^{TO}]$ is the expected timeout latency, and $R(\cdot)$ is the well-known TCP-friendly transmission rate derived in [23]. The detailed formula for computing $E[d_{ss}]$, $E[W_{ss}]$, $Q(\cdot)$, $E[Z^{TO}]$ and $R(\cdot)$ are omitted here due to space limitation — please refer to [11, 23].

B. Sensitivity Analysis for Short and Long TCP Flows

In this section, we analyze the transmission time for TCP flows of different sizes. Figure 2(a) gives in a log-log plot the average total latency computed from Equations (1) and (2) with $RTT = 0.1sec$, $RTO = 4RTT$ and $ITO = 3sec$, for a TCP flow of a fixed size FS for various p values.

We can observe that the average transmission time of short TCP flows are not very sensitive to loss when the loss rate is relatively small, but it increases drastically as loss rate becomes larger (when persistent congestion happens). For extremely long flows, the transmission time grows consistently with loss rate. Since it is hard, if not possible to compute the variance of transmission times, we use *ns* simulation [24] to study the coefficient of variation (C.O.V.)⁸ of transmission times for different flow sizes. In the simulation, a single TCP-Newreno connection transmits a file of size FS over a lossy link which randomly drops incoming packets with probability p . Packets never get dropped due to queue overflow. Once the file transfer finishes, a new connection immediately starts. We thus can measure the statistics of the transmission time. Figure 2(b) plots the C.O.V. against loss rate. We can observe an interesting trend: for small-size TCP flows, increasing the loss probability can lead to

⁷ For simplicity, we assume the receiver advertises a very large window value ($W_{max} \gg 1/p$).

⁸ C.O.V. is defined as the ratio between standard deviation and the mean of a random variable, smaller values imply less variability.

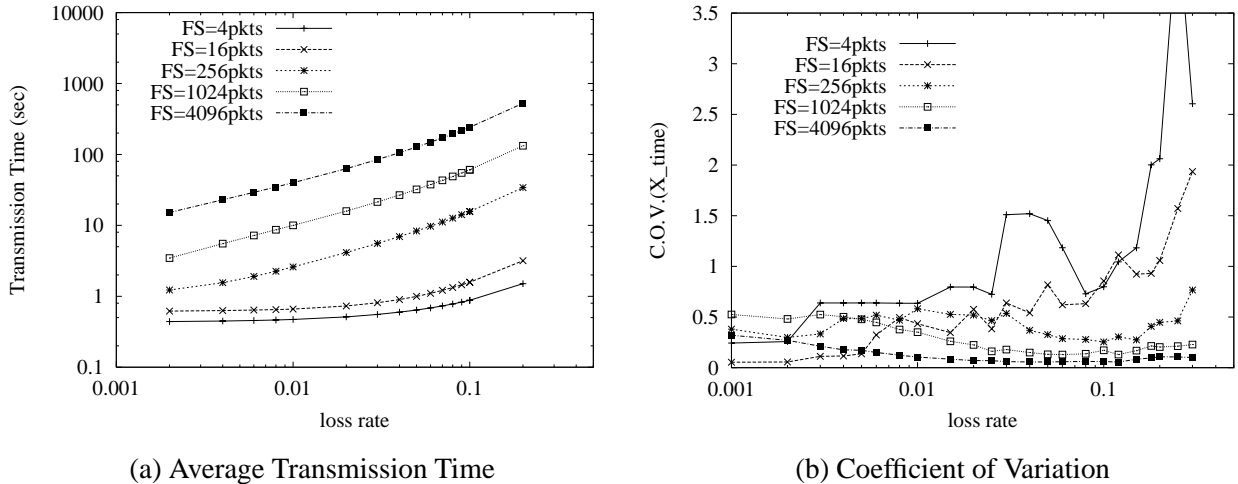


Fig. 2. Sensitivity Analysis of Transmission Time

increased variability, while for long TCP flows, large loss rate reduces the variability of transmission times. We can give a possible explanation to this interesting behavior:

- When loss rate is high, TCP congestion control is more likely to enter the exponential backoff phase, which can cause significantly high variability [25] in the transmission time of each individual packet of a flow.
- When loss rate is low, depending on when the packet loss happens, TCP can either transmit a significant amount of packets in slow-start phase or have to transmit them in the less aggressive congestion avoidance phase. This adds another dimension to variability, which affects the length of each congestion epoch and how many packets can be transmitted in each epoch (until another loss is detected).
- Since the first source of variability is on individual packets of a flow, the law of large numbers indicates that its impact is more significant on short flows. On the other hand, the second source of variability is more pronounced for long flows since most short flows finish their transmission in slow-start phase.

We thus conclude that reducing the loss probability is more critical to help short TCP flows experience less variations in transmission (response) time. Observe that the C.O.V. of transmission times is closely related to the fairness of the system—smaller values imply higher fairness. Such interesting behavior motivates us to give preferential treatment to short TCP flows.

C. Preferential Treatment to Short TCP Flows

In this section we use a simple simulation set up to illustrate that, assuming the arrival of requests does not depend on the service time, giving preferential treatment to short TCP flows can significantly enhance their transmission time, without degrading long flow performance.

For this purpose, we use the *ns* simulator [24] to set up the following scenario: 10 long (10000-packet) TCP-Newreno flows and 10 short (100-packet) TCP-Newreno flows are competing for bandwidth over a 1.25Mbps

link⁹. Flows start periodically and stay in the network until they finish transmission. We then vary the queue management policy at the bottleneck link and measure the instantaneous portion of bandwidth taken by each class of flows¹⁰ to show the effect of preferential treatment. We consider traditional Drop Tail queue, RED queue and our proposed RIO-PS (RIO with Preferential treatment to Short flows) queue (see Section III for details), and plot the results in Figure 3.

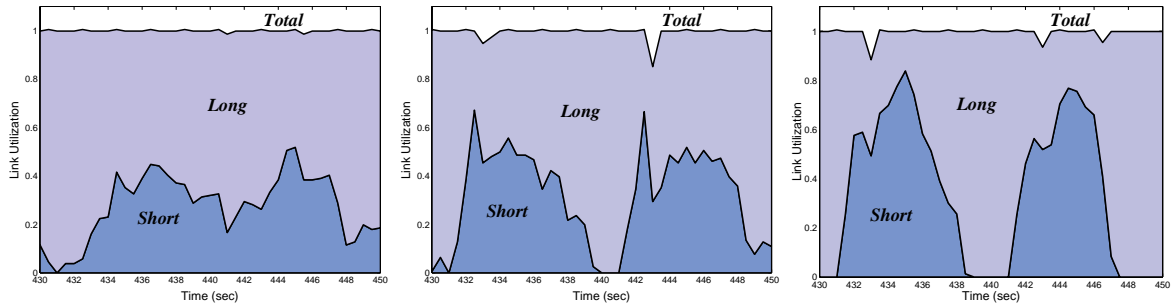


Fig. 3. Impact of Preferential Treatment— Link utilization under Drop Tail (left), RED (middle), and RIO-PS (right)

It is not surprising that a Drop Tail queue fails to give fair treatment to short TCP flows. Since most of the time short flows have a small congestion window compared to long flows, and also because a Drop Tail queue favors aggressive flows (with larger windows) [5], short flows can sometimes lose their fair share. A RED queue gives almost fair treatment to all flows, while our RIO-PS queue gives short flows more than their fair share.

As we can notice from the graphs, although when there are backlogs from both classes of flows, short flows under our scheme can temporarily grab (“steal”) more bandwidth from long flows, in the long run, their early completion returns an equal amount of resources (if not more, avoiding packet retransmissions) to long flows. In other words, since a long-lived TCP flow can adapt well to the network state and so can acquire about the same amount of resources (bandwidth) to transmit a certain file, giving preferential treatment to short flows does not sacrifice the long-term goodput¹¹ of long flows. In fact, this preferential treatment might even enhance the transmission of long flows since they operate in a more stable network environment (less disturbed by short flows) for longer periods. On the other hand, in a congested network, reducing the packet drops experienced by short flows can significantly enhance their response time and fairness among them, as implied by the analysis in the previous section. Table I gives the measured network goodput over the 500 seconds simulation period.

Table I also gives the measured goodput for a less loaded network with bottleneck link bandwidth of 1.5 Mbps. We observe that when load is low, our RIO-PS scheme (as well as RED) produces a slightly lower goodput (less than 1.5%) than that of DropTail. However, as load becomes higher, our scheme can even achieve slightly higher goodput than all other schemes¹².

⁹ Typically a short TCP flow has less than 20 packets to transmit. We only consider in this motivation section larger flow sizes to have a sharper contrast in our illustrations.

¹⁰ This instantaneous fraction of bandwidth is computed by the ratio of packets transmitted over link bandwidth over the last 0.5 second.

¹¹ Goodput is defined as the number of successfully delivered packets over the simulation period.

¹²Note that under Drop Tail, short flows are not protected well, so long flows achieve higher goodput than in other schemes.

We thus propose to differentiate between short and long TCP flows inside the network. To this end, we employ a Diffserv-like network architecture [26].

Link B/W	Flows	DropTail	RED	RIO-PS
1.25Mbps	All	153479	154269	154486
	Short	40973	49897	49945
	Long	112506	104372	104541
1.5Mbps	All	185650	184315	183154
	Short	43854	49990	49990
	Long	141796	134325	133164

TABLE I
NETWORK GOODPUT UNDER DIFFERENT SCHEMES

III. PROPOSED SCHEME: ARCHITECTURE AND MECHANISMS

In this section, we discuss the detailed implementation of our proposed scheme, including the network architecture we assume and the supporting mechanisms required to differentiate between short and long TCP flows.

A. The Architecture

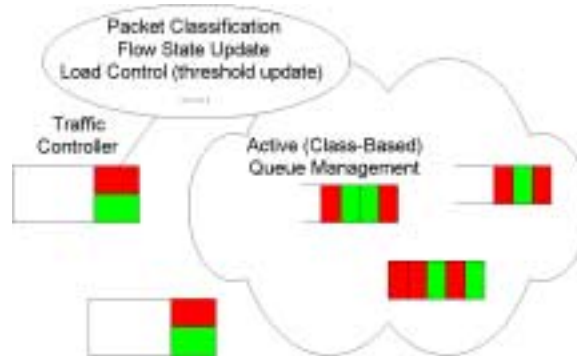


Fig. 4. Proposed Architecture

We assume a Diffserv-like scalable active domain management architecture where edge routers of an administrative domain perform all the per-flow information maintenance tasks, and the core routers only need to actively manage *per-class* flows. Figure 4 illustrates the operations a packet needs to go through in such architecture. The edge router tags (marks) individual packets as belonging to the class of short flows or the class of long flows. Such information is then used by core routers at which any class-based active queue management policy can be implemented.

B. Edge Router: Packet Classification and State Maintenance

We rely on the edge router to determine whether a packet is coming from a long or short flow. Accurate flow characterization can be very complicated due to the high variability of the user behavior in the current Internet.

Instead, we resort to a simple threshold-based approximation method. For each active flow, we simply maintain a counter that tracks how many packets have been observed so far. Once the counter exceeds a certain threshold (L_t), we consider the flow to be long, and classify the following packets as “Long” packets. Notice that under such classification method, the first L_t packets will be classified as “Short” packets. We discuss the pros and cons of such method in Section V.

The per-flow state information are softly maintained to detect the termination of a flow. More specifically, the flow hash table is updated periodically. Every T_u , a flow is considered terminated (and its entry removed) if no packets from that flow is observed in the last T_u time units.

The value of the threshold L_t can be static or dynamic. However, since the edge router already maintains enough per-flow state information, in our proposed scheme, we let the edge router adjust the threshold dynamically so as to balance the number of active short and active long flows. For such purpose, the edge router is configured with one parameter called Short-to-Long-Ratio (SLR) which is the ratio between the number of active short flows and that of long flows. The edge router then periodically (every T_c) performs additive increase additive decrease control over the threshold to achieve the target SLR. The pseudo-code for the control algorithm is shown in Appendix A. Other control policies are also possible, however, our simulation results show that the performance is not quite sensitive to SLR. Choosing appropriate values of T_u and T_c is an issue that needs further research (see Section V). In our simulation, we empirically choose $T_u = 1$ second and $T_c = 10$ seconds.

C. Core Router: Preferential Treatment to Short Flows

We require the core routers to give preferential treatment to short flows. There are many candidate queueing policies for such purpose. We choose the RIO (RED with In and Out) [4] queue since it has many advantages over other candidates. Figure 5 illustrates the early dropping/marking function of a RIO queue.

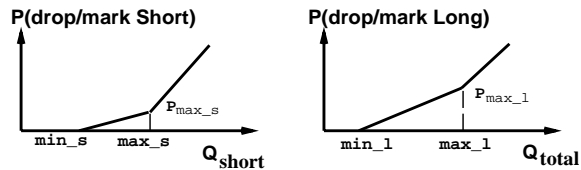


Fig. 5. RIO queue with Preferential treatment to Short Flows

Briefly speaking, the operation of the queue on In (Short) packets is not affected by the arrival of Out (Long) packets since the dropping/marking probability for Short packets is computed based on the average backlog of Short packets (Q_{short}) only. On the contrary, for the Long packets, the total average queue size (Q_{total}) is used to detect incipient congestion and thus flows that carry Long packets have to give up some resources (backoff) when there is persistent backlog from both classes of packets. The pseudo-code of the RIO algorithm can be found in [4], for completeness, we provide it in Appendix A. Note that we use the “gentle-” variant of RED dropping function as recommended in [27].

Clark and Feng [4] give a very nice discussion on the design features of RIO queue, and we only highlight few

of them here:

- Only a single FIFO queue is used for all the packets, thus packet reordering will not happen at a link even when packets from the same flow are classified differently.
- RIO inherits all the features of RED, which include protection of bursty flows, fairness within each class of traffic and detection of incipient congestion.
- RIO performs soft prioritization, thus does not lose the benefit of statistical multiplexing.

Henceforth, we will call such queue RIO-PS (RIO with Preferential treatment to Short flows). When persistent congestion happens, how much bandwidth share short flows can receive depends on the configured parameters for each class of traffic. In our simulation, we choose the RIO-PS parameters such that short TCP flows are guaranteed around 75% of the total link bandwidth in times of congestion.

IV. SIMULATION

We use the *ns* [24] simulator to study the performance of our proposed scheme and compare it to size-oblivious Drop Tail and RED schemes.

A. Simulation Setup

We assume network traffic is dominated by Web-like transactions. To this end, we adopt the Web traffic model developed in [28]. In this model, randomly selected clients initiate sessions which involve surfing several web pages of different sizes from randomly chosen websites. Each page may contain several objects, each of which requires a TCP connection for delivery (in other words, an HTTP 1.0 model is assumed). To request a page, the client sends a request packet to the server (simulating the GET message in a typical HTTP session), the server responds with an acknowledgment and then starts to transmit the web page requested by the client. The distributions of inter-page time, inter-object time, page and object sizes are given in Table II. The load is carefully tuned to be close to the bottleneck link capacity¹³. The topology used in our simulation is shown in Figure 6.

Name	inter-page	objs per page	inter-obj	obj size
Exp1 client 1	Exponential mean 9.5	Uniform min 2 max 7	Exponential mean 0.05	Bounded Pareto [4,200000] shape 1.2
Exp2 client 1	Exponential mean 9.5	Uniform min 2 max 7	Exponential mean 0.05	Bounded Pareto [4,500] shape 1.2
client 2	Exponential mean 192	Uniform min 1 max 3	Exponential mean 1.5	Bounded Pareto [400,200000] shape 1.2

TABLE II
WEB TRAFFIC CONFIGURATION

¹³ As revealed in [29], the advantage of RED is more pronounced in this load regime.

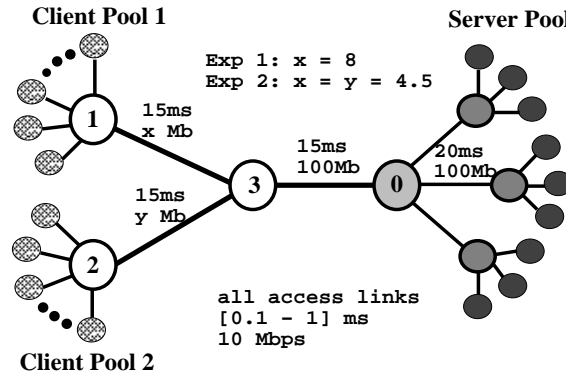


Fig. 6. RIO queue with Preferential treatment to Short Flows

Node 0 acts as the entry edge router of the bottleneck domain, and nodes 1,2,3 are core routers. Two sets of clients are connected to the server pool through individual bottleneck links (links 1-3 and 2-3, respectively). Other links are configured to be lossless. The buffer size of each bottleneck link and the configuration of its queue management policy are carefully tuned according to [29] so as to maximize power¹⁴. For our RIO-PS queue, we empirically choose its parameters so that short TCP flows experience less packet drops without starving long flows¹⁵. Recall that since the flow states are updated periodically at the edge routers, some terminated short flows may still be treated as “active”. Considering such measurement error due to lazy state update, we choose SLR to be 3 so that the actual number of short flows is approximately equal to that of long flows, however, we varied this ratio between 2 and 6 and didn’t observe significant difference. The detailed (default) configuration is listed in Table III. We have repeated each simulation with different random seeds and produced similar results¹⁶.

B. Experiment 1: Single Client set

In this experiment, there is only one set of clients involved (client pool 1). Therefore, the traffic seen at the core router 1 is the same as that at the edge router 0.

B.1 Study of the web traffic

We run the experiment for 4000 seconds simulation time. After some warm-up period (2000 seconds), we start to record the response time for each successfully downloaded web object. Figure 7(a) plots the averaged response time for different sized objects.

The results indicate that preferential treatment can cut the average response time for short and medium sized files significantly (25-30 %). In [2], Seddigh and Devetsikiotis observed that the default ITO value of current

¹⁴ Power is defined as the ratio between throughput and latency, so a high power implies low delay and high throughput.

¹⁵ In our simulations, we choose the RIO-PS parameters such that short TCP flows are guaranteed around 75% of the total link bandwidth. Even if short flows are given *absolute* priority, long flows may still pay very little penalty because of the heavy-tailedness of the workload, in which most of the bytes are due to the small percentage of long flows [18].

¹⁶ Also note that we only present the results of ECN capable routers here since ECN is recommended to be used with RED queues [30]. We also simulated the case when ECN was turned off and observed more pronounced performance gain of our proposed scheme over traditional RED, since RED with ECN turned on potentially enhances short flow performance [30].

Description	Value
Packet Size	500 bytes
Maximum Window	128 packets
TCP version	Newreno
TCP timeout Granularity	0.1 seconds
Initial Retransmission Timer	3.0 seconds
B/W delay product (BDP)	≈ 200 pkts (Exp1) ≈ 120 pkts (Exp2)
Bottleneck Buffer Size (B)	DropTail: $1.5 \times \text{BDP}$ RED/RIO-PS: $2.5 \times \text{BDP}$
Q. Parameters	$(min_{th}, max_{th}, P_{max}, w_q)$
RED	(0.15B, 0.5B, 1/10, 1/512)
RIO-PS short	(0.15B, 0.35B, 1/20, 1/512)
RIO-PS long	(0.15B, 0.5B, 1/10, 1/512)
RED & RIO-PS	ecn_on, wait_on, gentle_on
Edge Router	$SLR = 3, T_u = 1 \text{ sec}, T_c = 10 \text{ sec}$
Foreground Traffic	
(Src, Dest)	(Server Pool, Client Pool)
Long Connection Size	1000 packets
Short Connection Size	10 packets

TABLE III
NETWORK CONFIGURATION

implementations of TCP is too conservative and can thus significantly hurt the performance of short flows. They then propose to reduce the default ITO value from 3 seconds to 1 second. However, such aggressive value may lead to congestion collapse in some parts of the Internet where links are slow and have very large round trip propagation delay. Using aggressive retransmission timers at the ends of such links can cause unnecessary retransmissions, an undesirable phenomenon that TCP designers always wanted to avoid [31]. Nevertheless, we also run the simulation after setting the ITO to 1 second and plot the results in Figure 7(b). We observe that reducing the ITO value significantly reduces the performance gap between RED and our proposed scheme (now only 10-15%) for short connections (file size less than few dozens of packets). We still observe large performance improvements with RIO-PS for medium sized connections (15-25%).

Figure 8 plots the instantaneous queue size and average drop rate in the last 20 seconds for the case of 3 seconds ITO. Note to compare with DropTail, which does not support ECN, for RIO-PS and RED, we actually plot the sum of the forced drop (drop due to buffer overflow) and ECN mark rate. Also for RIO-PS we plot such measure for each class of flows.

We observe that our RIO-PS scheme reduces the overall mark/drop rate without the risk of congestion collapse (overloading) at the bottleneck link. A closer examination of the packet record reveals that the reduced dropping/mark rate comes from the fact that short connections rarely experience packet drops, and the number of ECN notifications sent to these connections are much smaller than their rivals. We argue that short connections are not responsible for controlling congestion since the timescale at which they are operating is much smaller than the effective control regime of a TCP congestion control algorithm. Thus, preferential treatment to short

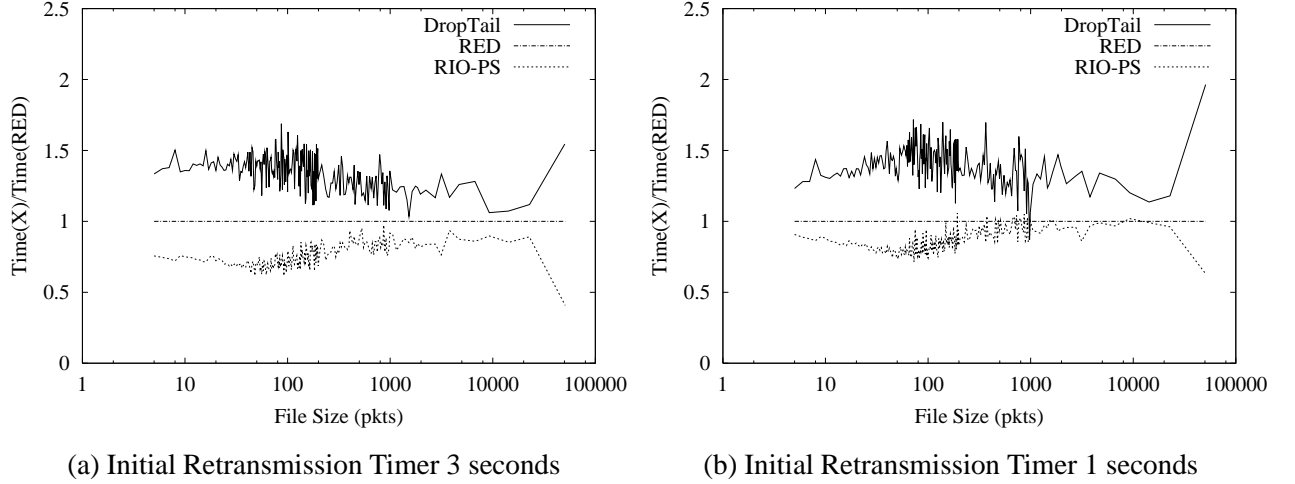


Fig. 7. Average response time relative to RED

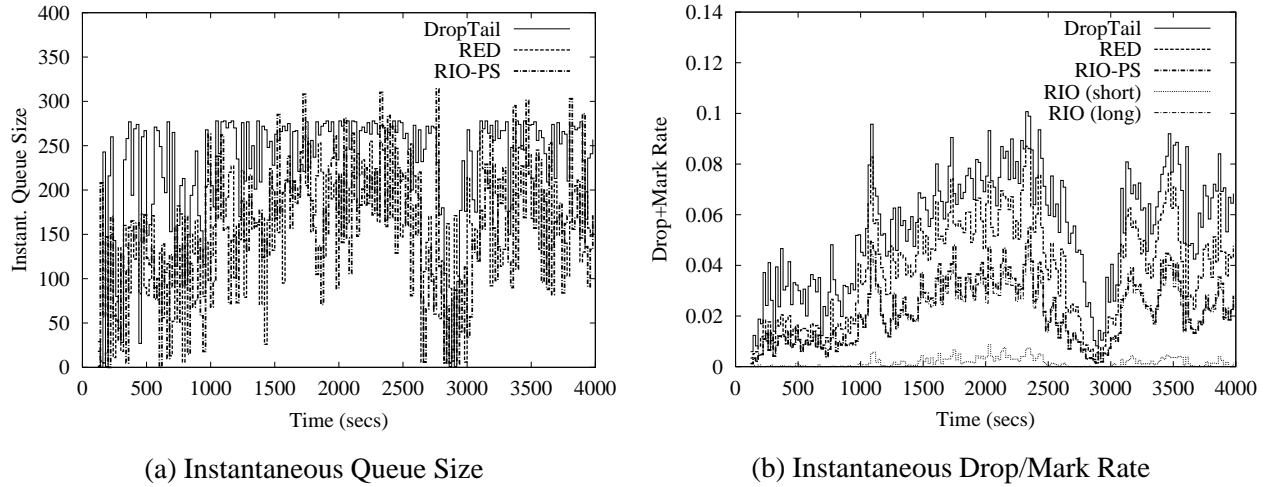


Fig. 8. Revealing the secret of better performance

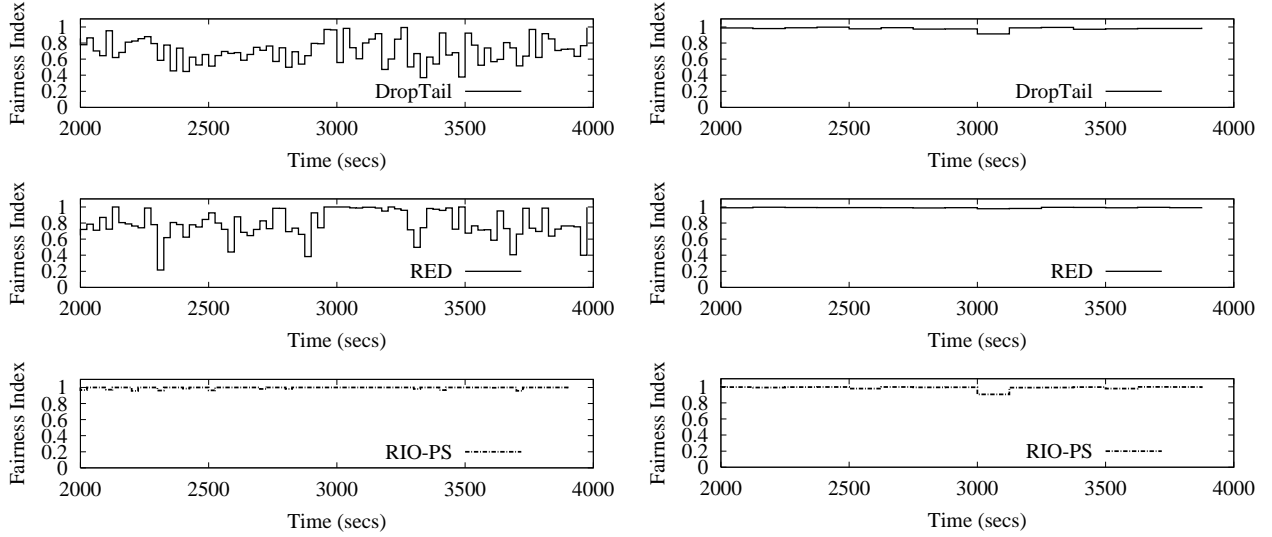
connections does not hurt, if not benefits the congestion state of the network.

B.2 Study of foreground traffic

To better illustrate the effect of the queue management policy on fairness of TCP connections, we periodically injected 10 short (every 25 seconds) and 10 long (every 125 seconds) foreground TCP connections and recorded the response time T_i for each connection i . The fairness index of response time, computed as $FI = \frac{(\sum_{i=1}^N T_i)^2}{N \sum_{i=1}^N T_i^2}$, is reported in Figure 9¹⁷. Figure 10 plots the transmission time for each individual connection and their ensemble average¹⁸.

¹⁷ Due to space limitation, the results for the case of $ITO = 1$ second are not presented here. The conclusion doesn't change, though.

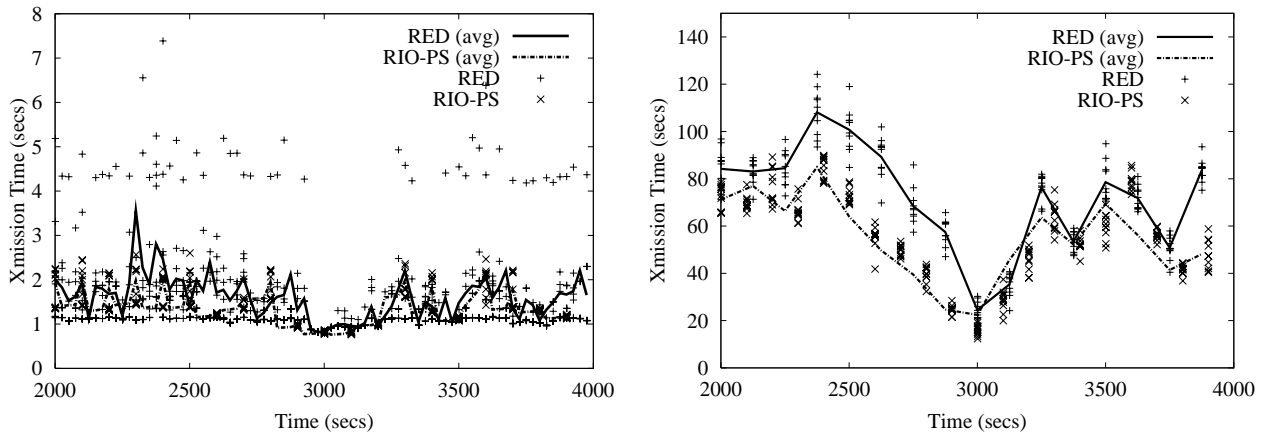
¹⁸ For the clarity of the graph, we don't plot the results of DropTail queues, whose performance (in terms of fairness and average transmission time) is much worse.



(a) Fairness Index of Short Connections

(b) Fairness Index of Long Connections

Fig. 9. Fairness of Transmission Time



(a) Transmission Time of Short Connections

(b) Transmission Time of Long Connections

Fig. 10. Transmission Time of Foreground Traffic

We observe that even with RED queues, many short flows experience at least one timeout, many of them (the “outliers” in Figure 10(a)) require waiting for a default ITO value (3 seconds) due to loss of the first packet¹⁹. On the contrary, our proposed preferential treatment scheme provides short TCP flows with much less drops and thus more fair service without hurting the performance of long connections.

¹⁹ Note that although ECN marks, instead of drops, packets in times of congestion, it does not mark the first (SYN) packet since the negotiation of ECN-capability at both ends has not finished at this moment. Therefore, initiating a TCP session when the network is congested still faces the risk of a long timeout if the SYN packet is dropped. Although *ns* does not generate control (SYN, SYN-ACK) packets, it similarly treats the first data packet, *i.e.* a RED with ECN scheme does not mark, rather it may drop, the first data packet in times of congestion.

C. Does fairness hurt goodput?

Table IV summarizes the overall goodput (averaged over 5 runs) of the network. As we can see, our proposed scheme does not hurt, if not (slightly) improves the overall goodput over Drop Tail.

Scheme	DropTail	RED	RIO-PS
Exp1 (ITO=3sec)	4207841	4264890	4255711
Exp1 (ITO=1sec)	4234309	4254291	4244158
Exp2 (ITO=3sec)	4718311	4730029	4723774

TABLE IV
NETWORK GOODPUT OVER THE LAST 2000 SECONDS

D. Experiment 2: Unbalanced Requests

The previous experiment assumes an ideal scenario for our proposed scheme in the sense that the portion of short connections seen by an edge router is the same as that by a core router. In a real ISP network, in which more complicated topologies are present, there may be unbalanced requests from different clients. We thus add another set of clients in our simulation model. We then unbalance the request load by designating client set 1 to only request relatively smaller objects (less than 500 packets) and client set 2 to request larger objects. Since Web sessions to different client sets take different routes (cf. Figure 6), the load control at the edge router is now not as effective as in the previous case. We consider such case as the worst-case scenario for our proposed scheme. Figure 11 plots performance measures in such scenario for 10-packet foreground TCP connections to client set 1 and 1000-packet foreground TCP connections to client set 2.

When the route is dominated by one class of flows (short or long), preferential treatment in core routers no longer has significant impact on the flow performance. In other words, in such extreme case, our proposed scheme reduces to traditional unclassified traffic plus RED queue policy. However, giving preferential treatment to the initial few packets still helps to reduce the chance of retransmitting the first packet and thus helps most short connections finish quickly, as shown in Figure 11.

V. DISCUSSION

A. Comments on the Simulation Model

Our simulation is performed using the “Dumbbell and Dancehall” (one-way traffic) model [32], and all TCP connections have similar end-to-end propagation delays, thus such model may not represent the common topology seen by Internet users. However, such assumption is only for the sake of easy comparison of individual TCP performance. When there is also reverse traffic present, our proposed scheme will even have more superior performance over other traditional policies. This is because the setup of a web session (and most other TCP sessions) requires the client side to exchange a short sequence of control packets with the server. Such a short exchange will be better protected under our proposed scheme, especially in the presence of heavy traffic in the same

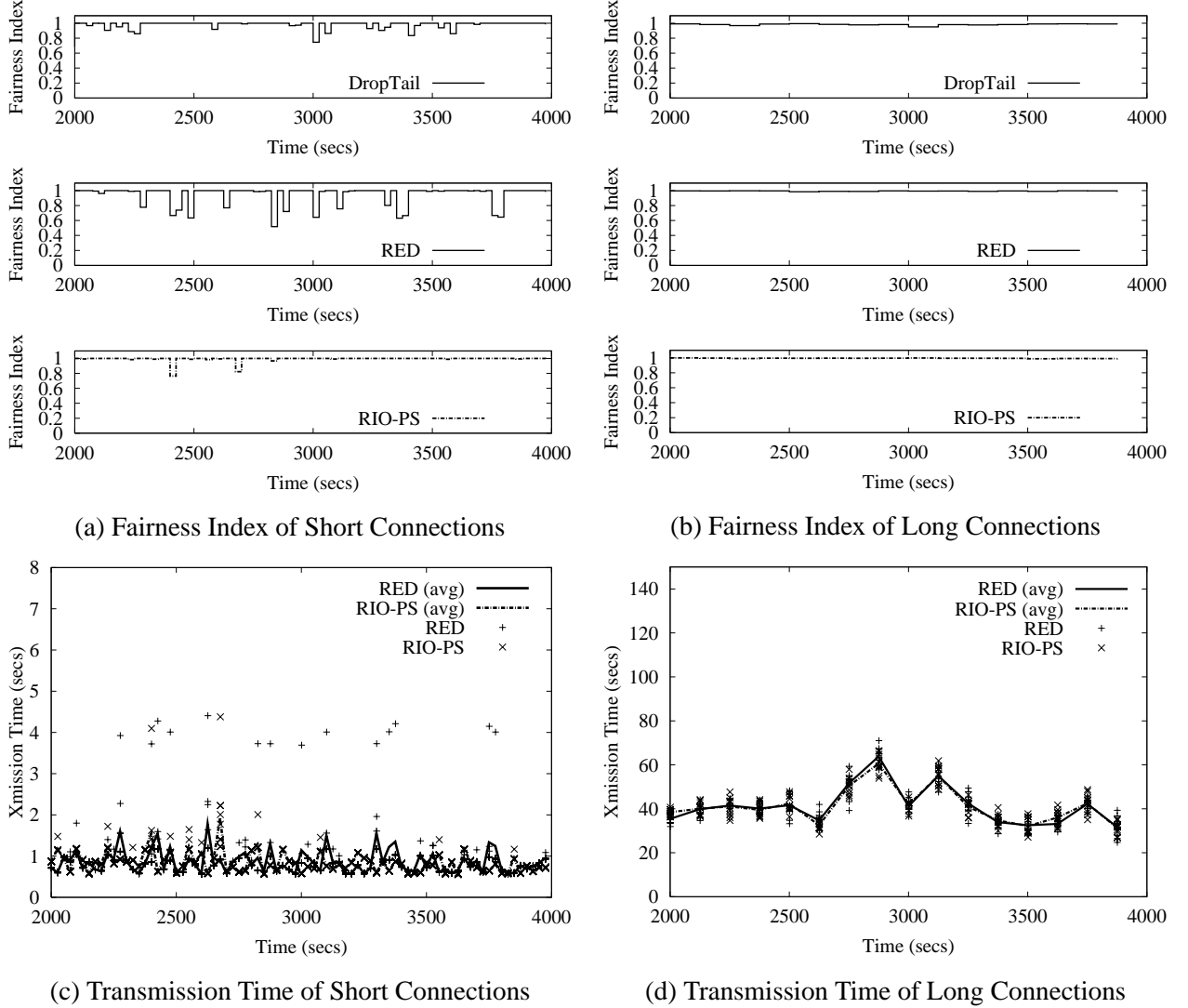


Fig. 11. Foreground Traffic performance comparison

direction. More complicated topology models may have a negative impact on the performance of our scheme due to unbalanced requests. However, as we pointed out in Section IV-D, our scheme can still work at least as well as traditional policies even in extremely unbalanced cases.

B. Deployment Issues

Our proposed scheme requires edge devices be able to perform per-flow state maintenance and per-packet processing. Fortunately, it has been shown in [33] that providing such functionality does not have a significant impact on the end-to-end performance. The other issue is on whether such scheme is incrementally deployable. Notice that our objective is to protect short delay-sensitive flows, and our success relies on how well the classification (isolation) is done at the edge (bottleneck) router. The scheme does not require the queue mechanisms be implemented at each router. For example, the advanced edge device can be placed in front of a busy web server

cluster, and RIO-PS only implemented at the busy bottleneck links to achieve the desired performance objective.

C. Flow Classification

We use a threshold-based classification method since the intermediate (edge) node can not predict in advance whether a flow is long or short. Such method thus mistakenly classifies the first few packets of a long flow as if they came from a short flow. However, such “mistake” may help to enhance performance. This is because the first few packets of a long TCP flow is under exactly the same control method as packets from a short flow. Thus, they are more vulnerable to packet losses and deserve to be treated with high preference. Such “mistake” also makes the system more fair to all TCP connections.

D. Controller Design

Our proposed scheme involves controller design issues at different places and timescales. The issue of edge load control is certainly a topic for further research, although our preliminary results indicate that the performance is not very sensitive to the target load ratio of active short to active long flows (the value of SLR at the edge). The “actual” SLR depends on the values of T_c and T_u , which determine how often the classification threshold L_t and active flow table are updated, respectively. Smaller values of T_c and T_u may increase accuracy at the expense of increased overhead.

AQM policies other than RIO may also be used in our scheme. For example, RED with PI control [34] can be used for long flows since the controller was designed specially for long-lived connections. We are currently studying the performance of such controller in the presence of both short and long TCP flows.

VI. CONCLUSIONS AND FUTURE WORK

TCP traffic constitutes the majority of the bytes flowing over the Internet today. In this paper, we present a Diffserv-like architecture where edge routers implement a new TCP service by classifying TCP flows based on their size (lifetime). Core routers implement a simple RIO (weighted RED) policy to give preferential treatment to packets from short TCP flows. This has several advantages over size-oblivious Drop Tail or RED schemes: (1) the performance of the majority of TCP flows (the short transfers or mice) is significantly improved in terms of response time and fairness; (2) because short TCP flows can be rapidly served, the performance of the few TCP long flows (the large transfers or elephants) is also enhanced or at worst minimally affected; (3) the overall goodput of the system is thus improved or at least stays almost the same (a work-conservation feature); and (4) the proposed architecture is extremely flexible in that the functionality that defines the new TCP service can be largely tuned only at the edge routers.

The size-aware management of traffic has shown considerable promise. We have shown how differentiation of TCP flows based on size can be effectively used to isolate flows that mostly operate in different regimes: short TCP flows in slow-start and long TCP flows in congestion-avoidance. At the network layer, Shaikh *et al.* [35] investigated load-sensitive routing of only long-lived IP flows in order to improve routing stability and reduce

overhead. However, they have not considered the effect on congestion-controlled sources (such as TCP and TCP-friendly applications). We are currently investigating an approach that integrates size-aware traffic management at both the network and transport layers.

REFERENCES

- [1] R. Khare and I. Jacobs, "W3C Recommendations Reduce 'World Wide Wait'," <http://www.w3.org/Protocols/NL-PerfNote.html>.
- [2] N. Seddigh and M. Devetsikiotis, "Studies of TCP's Retransmission Timeout Mechanism," in *Proc. ICC 2001*, Helsinki, Finland, June 2001.
- [3] Y. Bernet and et al., "A Framework for Differentiated Services," Internet Draft (draft-ietf-diffserv-framework-02.txt), February 1999.
- [4] D.D. Clark and W. Fang, "Explicit Allocation of Best-Effort Packet Delivery Service," *IEEE/ACM Transactions on Networking*, vol. 6(4), pp. 362–373, August 1998.
- [5] S. Floyd and V. Jacobson, "Random Early Detection Gateways for Congestion Avoidance," *IEEE/ACM Transactions on Networking*, vol. 1(4), pp. 397–413, August 1993.
- [6] S. Floyd and V. Jacobson, "Link-sharing and Resource Management Models for Packet Networks," *IEEE/ACM Transactions on Networking*, vol. 3, no. 4, pp. 365–386, August 1995.
- [7] I. Matta and L. Guo, "Differentiated Predictive Fair Service for TCP Flows," in *Proc. ICNP'2000*, Osaka, Japan, November 2000.
- [8] T.S.E. Ng, D.C. Stephens, I. Stoica, and H. Zhang, "Supporting Best-Effort Traffic with Fair Service Curve," in *Proc. IEEE GLOBECOM'99*, Rio de Janeiro, Brazil, December 1999.
- [9] B. Nandy, N. Seddigh, A. Chapman, and J. Hadi Salim, "A Connectionless Approach to Providing QoS in IP Networks," in *Proc. IFIP Eighth Conference on High Performance Networking (HPN'98)*, Vienna, September 1998.
- [10] V. Firoiu and X. Zhang, "Best Effort Differentiated Services: Trade-off Service Differentiation for Elastic Applications," Tech. Rep., Nortel Networks, September 2000.
- [11] N. Cardwell, S. Savage, and T. Anderson, "Modeling TCP Latency," in *Proc. IEEE INFOCOM 2000*, Tel-Aviv, Israel, March 2000.
- [12] B. Sikdar, S. Kalyanaraman, and K.S. Vastola, "An Integrated Model for the Latency and Steady State Throughput of TCP Connections," in *Proc. IFIP Symposium on Advanced Performance Modeling*, Orlando, FL, November 2000.
- [13] M. Allman, S. Floyd, and C. Partridge, "Increasing TCP's Initial Window," Internet RFC 2414, September 1998.
- [14] V. Padmanabhan and R. Katz, "TCP Fast Start: a Technique for Speeding Up Web Transfers," in *Proc. IEEE Globecom'98 Internet Mini-Conference*, November 1998.
- [15] Y. Zhang, L. Qiu, and S. Keshav, "Speeding Up Short Data Transfers: Theory, Architecture Support, and Simulation Results," in *Proc. NOSSDAV 2000*, Chapel Hill, NC, June 2000.
- [16] J. Heidemann, "Performance Interactions Between P-HTTP and TCP Implementations," *ACM Computer Communication Review*, vol. 27(2), pp. 65–73, April 1997.
- [17] M. Crovella, M. Harchol-Balter, and C. Murta, "Task Assignment in a Distributed System: Improving Performance by Unbalancing Load," in *Proc. ACM Sigmetrics '98 Poster Session*, Madison, WI., June 1998.
- [18] M. Crovella, R. Frangioso, and M. Harchol-Balter, "Connection Scheduling in Web Servers," in *Proc. the 1999 USENIX Symposium on Internet Technologies and Systems (USITS '99)*, Boulder, Colorado, October 1999.
- [19] G.L. Monoco, F. Azeem, and S. Kalyanaraman, "TCP-Friendly Marking for Scalable Best-Effort Services on the Internet," *ACM Computer Communication Review*, to appear 2001.
- [20] T. Bonald and L. Masoulié, "Impact of Fairness on Internet Performance," in *Proc. ACM SIGMETRICS 2001*, Boston, MA, June 2001.
- [21] W. Stevens, *TCP/IP Illustrated, Vol. 1: The Protocols*, Addison-Wesley, 1997.
- [22] R. Braden, "Requirements for Internet Hosts – Communication Layers," Internet RFC 1122, October 1989.
- [23] J. Padhye, V. Firoiu, D. Towsley, and J. Kurose, "Modeling TCP Throughput: A Simple Model and its Empirical Validation," in *Proc. ACM SIGCOMM'98*, September 1998.
- [24] E. Amir et al., "UCB/LBNL/VINT Network Simulator - ns (version 2)," Available at <http://www.isi.edu/nsnam/ns/>.
- [25] L. Guo, M. Crovella, and I. Matta, "How does TCP Generate Pseudo-Self-Similarity," Tech. Rep. 2000-017, Dept. of Computer Science, Boston University, Boston, MA, July 2000, Submitted for Publication.
- [26] S. Black et al., "An Architecture for Differentiated Services," Internet RFC 2475, December 1998.
- [27] S. Floyd, "Recommendation on using the 'gentle.' variant of red," <http://www.aciri.org/floyd/red/gentle.html>, March 2000.
- [28] A. Feldmann, A.C. Gilbert, P. Huang, and W. Willinger, "Dynamics of IP Traffic: A Study of the Role of Variability and the Impact of Control," in *Proceeding of ACM SIGCOMM'99*, Boston, Mass, September 1999.
- [29] M. Christiansen, K. Jeffay, D. Ott, and F.D. Smith, "Tuning RED for Web Traffic," in *Proc. ACM SIGCOMM 2000*, Stockholm, Sweden, Aug.-Sep. 2000.
- [30] S. Floyd, "TCP and Explicit Congestion Notification," *ACM Computer Communication Review*, vol. 24(5), October 1994.
- [31] V. Jacobson, "Congestion Avoidance and Control," in *Proc. ACM SIGCOMM'88*, Stanford, CA, August 1988.
- [32] K. Nichols, "Dumbbells, Dancehalls, and Beyond: Understanding QoS Prior to Deployment," in *Proceedings of the Internet Bandwidth Management Summit (iBAND3)*, San Francisco, CA, October 1999.
- [33] R. Guérin, L. Li, S. Nadas, P. Pan, and V. Peris, "The Cost of QoS Support in Edge Devices: An Experimental Study," in *Proceedings of IEEE INFOCOM'99*, New York, NY, March 1999.
- [34] C.V. Hollot, V. Misra, D. Towsley, and W. Gong, "On Designing Improved Controllers for AQM Routers Supporting TCP Flows," in *Proc. IEEE INFOCOM 2001*, Anchorage, AL, April 2001.

- [35] A. Shaikh, J. Rexford, and K.G. Shin, "Load-Sensitive Routing of Long-Lived IP Flows," in *Proceedings of SIGCOMM'99*, Boston, MA, September 1999.

APPENDIX

I. PSEUDO-CODE FOR TCP DIFFERENTIATION MECHANISMS

A. Edge Router load control module

```

1. Upon update timer expiration (every  $T_c$  units):
2.   Scan flow table,
3.   Count the number of long (short) flows  $cnt_l$  ( $cnt_s$ ).
4.   If  $cnt_s > cnt_l \times SLR$ 
5.      $L_t \leftarrow L_t - 1$ 
6.   Else
7.      $L_t \leftarrow L_t + 1$ 

```

B. Core RIO-PS queue dropping function

Denote by q_{short} and q_{total} the instantaneous number of short packets in queue and the total queue size, respectively.

```

1. For each arrived packet  $P$ 
2.   If  $P \in Short$ 
3.      $Q_{short} \leftarrow w_{short} \cdot q_{short} + (1 - w_{short}) \cdot Q_{short}$ 
4.      $Q_{total} \leftarrow w_{total} \cdot q_{total} + (1 - w_{total}) \cdot Q_{total}$ 
5.   If  $P \in Short$ 
6.     If  $min_{short} < Q_{short} < max_{short}$ 
7.        $p_s \leftarrow calculate\_p(Q_{short}, min_{short}, max_{short}, P_{max})$ 
8.       with probability  $p_s$  drop  $P$ .
9.     Else If  $Q_{short} > max_{short}$ 
10.      drop  $P$ .
11.   Else
12.     If  $min_{long} < Q_{total} < max_{long}$ 
13.        $p_l \leftarrow calculate\_p(Q_{total}, min_{long}, max_{long}, P_{max})$ 
14.       with probability  $p_l$  drop  $P$ .
15.     Else If  $Q_{total} > max_{long}$ 
16.      drop  $P$ .

```