

# Tuning RED for Web Traffic

Mikkel Christiansen, Kevin Jeffay, David Ott, and F. Donelson Smith

**Abstract**—We study the effects of RED on the performance of Web browsing with a novel aspect of our work being the use of a user-centric measure of performance: response time for HTTP request–response pairs. We empirically evaluate RED across a range of parameter settings and offered loads. Our results show that: 1) contrary to expectations, compared to a FIFO queue, RED has a minimal effect on HTTP response times for offered loads up to 90% of link capacity; 2) response times at loads in this range are not substantially affected by RED parameters; 3) between 90% and 100% load, RED can be carefully tuned to yield performance somewhat superior to FIFO, however, response times are quite sensitive to the actual RED parameter values selected; and 4) in such heavily congested networks, RED parameters that provide the best link utilization produce poorer response times. We conclude that for links carrying only Web traffic, RED queue management appears to provide no clear advantage over tail-drop FIFO for end-user response times.

**Index Terms**—Active queue management, congestion control, experimental methodology, HTTP, random early detection (RED).

## I. INTRODUCTION

A RECENT IETF publication strongly recommended the widespread deployment of active queue management technology in routers to improve the performance of today's Internet [4]. Active queue management refers to the practice of manipulating the queue at an outbound interface in a router to bias the performance of flows that transit the router. The goals of active queue management are to 1) reduce the average length of queues in routers and thereby decrease the end-to-end delay experienced by packets; and 2) ensure that network resources are used more efficiently by reducing the packet loss that occurs when queues overflow.

The recommended active queue management to be deployed is *random early detection*, better known as RED [13]. Under RED, a router will probabilistically drop an arriving packet even though the queue for the appropriate outbound interface is not full. The motivation for this “early” drop comes from the fact that packet loss is the primary indicator of congestion for a TCP connection. By dropping packets before a router's queue fills, the TCP connections sharing the queue will reduce their transmission rates and (ideally) ensure the queue does not overflow. The claim (borne out by significant empirical data) is that dropping packets prior to the overflow of the queue will reduce the

overall rate of packet loss. Given that TCP traffic dominates on Internet backbones [25], RED and other forms of early congestion notification have the potential to improve overall network performance as well as that seen by individual TCP connections. In this work, we test this claim and explore the impact of RED on the performance of the most dominant subset of TCP connections on the Internet today: Web traffic. In particular, we are interested in measuring the effect of RED on a user-centric measure of performance: the response time for an HTTP 1.0 request. Although the performance of RED and other early congestion notification mechanisms continue to be the subject of much study, the evaluation metrics have largely been network-centric measures such as network link utilization or aggregate TCP throughput. Moreover, as argued in Section II, most of these evaluation studies focused on simulations of long-lived TCP connections such as (huge) file transfers. In contrast, measurement studies have shown that the majority of TCP connections are HTTP connections<sup>1</sup> and that many of these connections are quite short lived, often on the order of a few TCP segments. More importantly, given that the performance of the Internet is becoming synonymous with the performance of the Web, understanding the impact of router forwarding behaviors on user-visible performance measures is an important (and largely ignored) aspect of the evaluation of any congestion control proposal.

At a high level, we seek to compare the performance of HTTP request–response pairs under RED and more traditional tail-drop<sup>2</sup> first-in-first-out (FIFO) queueing. Unfortunately, measuring the performance of HTTP under RED is a complex problem. First, as described in more detail in Section III, RED is a general mechanism that is controlled by (at least) five separate control parameters. There exist rules-of-thumb for assigning values to most parameters [15], but little is known about how (or if) one can optimize RED performance for a given traffic class. Second, even if optimal RED parameter settings were known, generating or simulating HTTP behaviors in a meaningful way is problematic. There are few models of HTTP traffic and it is likely the case that Web traffic dynamics (e.g., the mix between HTTP 1.0 and 1.1 protocols) are evolving faster than our current ability to measure and model the traffic.

Our general approach is to conduct a “live simulation” of Web browsing in a laboratory environment. By live simulation, we mean that we simulate a large collection of users browsing the Web at a set of sites distributed throughout the continental U.S. The HTTP traffic generated by the simulated users will traverse a laboratory network with routers that support both RED and traditional tail-drop FIFO queueing. A number of instances of the

Manuscript received July 20, 2000; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor C. Diot. This work was supported by the National Science Foundation under Grants CDA-9624662 and ITR 00-82870, and by Cisco Systems, the Intel Corporation, and the North Carolina Networking Initiative.

M. Christiansen was with the Department of Computer Science, University of North Carolina, Chapel Hill, NC 27599-3175 USA. He is now with the Department of Computer Science, Aalborg University, DK-9220 Aalborg, Denmark.

K. Jeffay, D. Ott, and F. D. Smith are with the Department of Computer Science, University of North Carolina, Chapel Hill, NC 27599-3175 USA.

Publisher Item Identifier S 1063-6692(01)04724-0.

<sup>1</sup>For example, recent measurements on the MCI backbone show that about 95% of the bytes transmitted across the network are carried by TCP and of these, 50%–70% are HTTP messages [25].

<sup>2</sup>Tail-drop refers to the action of dropping a packet that was intended to be added at the tail of a full FIFO queue.

user-browsing simulation program are run to generate a configurable offered load on a bottleneck network link. The user HTTP requests will be delivered to a set of servers that will respond with responses of the appropriate sizes. Both request and response packets are artificially delayed to simulate the round-trip times (RTT) experienced when communicating with machines distributed across the U.S. This is done to ensure our end-to-end response-time measurements reflect the full range of effects of TCP congestion control and retransmissions experienced by real users. When the responses are delivered back to the users, we record the elapsed time for each simulated HTTP request–response pair.

This experimental setup provides a basis for comparing the effect of RED versus FIFO queueing on the response time for HTTP requests. We performed a series of experiments to empirically determine the FIFO queue length and combination of RED parameter settings that result in the best performance for our network and our simulation of Web traffic. From our experiments, we observe the following.

- Contrary to expectations, when compared to a (properly configured) tail-drop FIFO queue, RED has a minimal effect on HTTP response times for offered loads up to 90% of link capacity.
- Response times for loads in this range are not substantially affected by values of RED parameters.
- Between loads of 90% to 100% of link capacity, RED can be carefully tuned to yield performance somewhat superior to FIFO. However, response times are quite sensitive to the actual RED parameter values selected. In our experiments recommended parameter settings resulted in poorer performance than FIFO. Worse, the “optimal” settings that resulted in the best RED performance were nonobvious and arrived at only through exhaustive trial-and-error experimentation.
- For loads of 90% to 100% of link capacity where RED has the potential to provide better performance, performance becomes a subjective measure. For loads in this range there exists a tradeoff between improving response times of short-lived connections and improving response times of long-lived connections. Both cannot be optimized simultaneously.
- In such heavily congested networks, there exists a tradeoff between network utilization and HTTP transaction response times. RED parameters values that provide the best link utilization produce poorer response times.

We have considered only HTTP traffic in our experiments and hence our results are best interpreted as representing a worst-case scenario for RED performance on real Internet links that carry a mix of HTTP and other traffic classes. Nonetheless, we conclude that for links carrying only Web traffic, RED appears to provide no clear advantage over tail-drop FIFO for end users whose primary metric of satisfaction is response time. Moreover, given the lack of engineering practice to guide the setting of RED parameter values, and our demonstration that “reasonable” but nonetheless suboptimal RED parameters values can result in poorer performance than FIFO queueing, without further analysis it is possible that widespread RED deployment may not provide the expected benefits.

The remainder of this paper is organized as follows. Section II provides a more in-depth introduction to RED and reviews the literature in the performance evaluation of RED and related active queue management schemes. Section III describes our experimental methods and the design and calibration of our experiments. Section IV presents the performance of our simulated Web browsing sessions under FIFO queueing, and Section V presents results for RED queueing. Section VI provides a preliminary analysis of the observed response time distributions under RED. Section VII compares the RED and FIFO results. We conclude in Section VIII with a discussion of the results, the limitations of our experiments and results, and some comments on future work.

## II. BACKGROUND AND RELATED WORK

RED is an active queue management mechanism that is intended to address some of the shortcomings of standard tail-drop FIFO queue management [4]. In a FIFO queue it is possible for “lock-out” to occur, a condition in which a small subset of the flows sharing the link can monopolize the queue during periods of congestion. Flows generating packets at a high rate can fill up the queue such that packets from flows generating packets at substantially lower rates have a higher probability of arriving at the queue when it is full and being discarded. A second problem with a FIFO queue is that latency is increased for all flows when the queue is constantly full. Simply making the queue shorter will decrease the latency but negates the possibility of accommodating brief bursts of traffic without dropping packets unnecessarily. RED addresses both the lock-out problem by using a random factor in selecting which packets to drop and the “full queue” problem by dropping packets early, before the queue fills.

The RED algorithm uses a weighted average of the total queue length to determine when to drop packets. When a packet arrives at the queue, if the weighted average queue length is less than a minimum threshold value,  $\min_{th}$ , no drop action will be taken and the packet will simply be enqueued. If the average is greater than  $\min_{th}$  but less than a maximum threshold,  $\max_{th}$ , an *early drop* test will be performed as described below. An average queue length in the range between the thresholds indicates some congestion has begun and flows should be notified via packet drops. If the average is greater than the maximum threshold value, a *forced drop* operation will occur. An average queue length in this range indicates persistent congestion and packets must be dropped to avoid a persistently full queue. (The forced drop is also used when the queue is full but the average queue length is still below the maximum threshold.) Note that by using a weighted average, RED avoids overreaction to bursts and instead reacts to longer term trends. Furthermore, because the thresholds are compared to the weighted average (with a typical weighting factor  $w_q$  of 1/512), it is possible that no forced drops will take place even when the instantaneous queue length is quite large. For example, Fig. 1 illustrates the queue length dynamics in a RED router used in our experiments. For the experiment illustrated in Fig. 1, forced drops would occur only in the one short interval near the beginning when the weighted average reaches the maximum threshold.

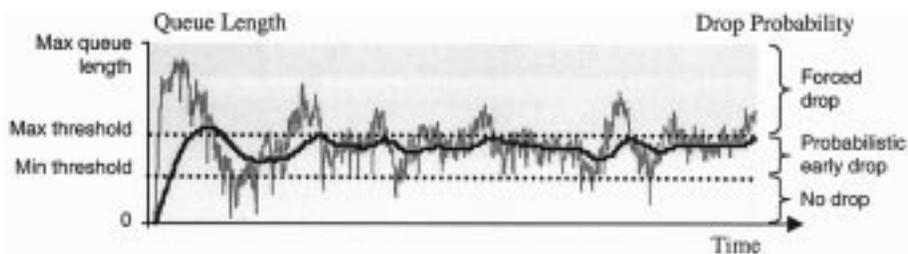


Fig. 1. Illustration of the desired behavior of a packet queue at an outbound interface in a router. The gray line indicates instantaneous queue length; the black line indicates the weighted average queue length.

TABLE I  
RED CONTROL PARAMETERS

$qlen$	The maximum number of packets that can be enqueued.
$min_{th}$	Queue length threshold for triggering probabilistic drops.
$max_{th}$	Queue length threshold for triggering forced drops.
$w_q$	Weighting factor for the average queue length computation.
$max_p$	The maximum probability of performing an early drop.

The early drop action in the RED algorithm probabilistically drops the incoming packet when the weighted average queue length is between the  $min_{th}$  and  $max_{th}$  thresholds. In contrast, the forced drop action in the RED algorithm is guaranteed to drop the incoming packet. In the case of early drops, the probability that the packet will be dropped is dependent on several other parameters of the algorithm. An initial drop probability  $P_b = max_p \cdot (avg - min_{th}) / (max_{th} - min_{th})$  is computed, where  $max_p$  is the maximum drop probability (an additional control parameter) and  $avg$  is the weighted average queue length. The actual drop probability is a function of the initial probability and a count of the number of packets enqueued since the last packet was dropped:  $P_a = P_b / (1 - count \times P_b)$ . Note that given a weighted average queue size, the impact of  $min_{th}$  is dependent on both  $max_p$  and  $max_{th}$ . This means that one may find a value for  $min_{th}$  that results in good performance, but it may only be in combination with certain values of  $max_p$  and  $max_{th}$ . In principle, this is the case for all the parameters. The main control parameters for RED are summarized in Table I.

The design of RED is such that during the drop phases of the algorithm, high-bandwidth flows will have a higher number of packets dropped since their packets arrive at a higher rate than lower bandwidth flows (and thus are more likely to be dropped in an early drop action). However, all flows experience the same loss rate under RED. By using probabilistic drops, RED maintains a shorter average queue length, avoiding lockout and repeated penalization of the same flows when a burst of packets arrives.

The original RED paper [13] presented analysis and several simulations to show the results of RED usage and develop insights into the effects different RED parameters have on performance. They arrived at suggested guidelines for useful ranges of parameter values and explanations of the considerations that

would influence tuning parameters to achieve desired results for particular traffic characteristics. Subsequent analysis by RED's designers and others led to the current guidelines [15] that are discussed later in this paper.

One of the earliest experiments with RED was reported in [26] and gives the results of live testing with a RED implementation in a router ahead of a bottleneck DS3 link in a transcontinental network. These tests were conducted with a small number of continuously sending high-bandwidth TCP connections. Total throughput of the TCP connections was the primary measure of performance and delays were not measured. The results showed that, in general, RED achieved better throughput and better link utilization for multiple connections than comparable tail-drop FIFO. RED was also effective in preventing congestion collapse when the TCP windows were configured to exceed the storage capacity of the network. A very important result showed that the interface queue (buffer) size is a critical parameter even with RED and should be 1–2 times the bandwidth-delay product at a bottleneck link.

A number of research efforts have focused on possible shortcomings of the algorithms in RED and have proposed modifications and alternatives, among them BLUE [12], Stabilized RED (SRED) [23], Adaptive RED [11], FRED (Flow Random Early Drop) [17], and BRED (Balanced RED) [2]. We do not comment here on the contributions and merits of these proposals except to note any analysis or simulations that examine the behavior and performance of "classic" RED. For example, in [17] simulations are used to demonstrate situations in which RED does not provide protection from nonadaptive flows, and situations in which RED does not promote fair sharing of link bandwidth between TCP connections with long RTT or small windows, and other competing flows.

In [12], there are a suite of results from ns simulations of RED with *explicit congestion notification* (ECN) [14] enabled in both routers and end-system TCP implementations. The simulations focused primarily on the effects of the parameter  $w_q$  used to smooth measurements of the average queue size. Interestingly, some of these simulations use a large number of sources (1000–4000) that generate traffic with Pareto on/off periods and might provide clues to behavior in Web-like traffic. Unfortunately, because all the simulations use ECN marking instead of packet drops, and end-to-end delays are not considered, the results are not directly comparable to our work on packet-drop RED. Feng *et al.* presents ns simulations of RED with packet drops in situations where a moderate number (32 or 64) of continuously sending TCP connections share a link [11]. Here the

maximum drop probability  $\max_p$ , was varied to see its effect on loss rates and average queue length. The results show that the “best” value for  $\max_p$  is dependent on the number of connections and, for any setting, the drop rate is not significantly different from that of a tail-drop FIFO queue. The argument is also made that the effectiveness of RED decreases as the number of connections sharing the queue increases. This is because a small number of connections actually receive and act on RED-induced congestion indications.

Results reported in [23] for simulations of RED with persistent (continuously sending) TCP connections (ranging from 10–1000 connections) showed that router queue lengths (measured in the total buffer space consumed) were at or below the minimum threshold for a small number of connections and stabilized around the maximum threshold for a large number of connections. Simulations were also conducted with more “realistic” traffic by using a large number of TCP connections (2000–3500) to transfer random size files with a size distribution derived from measurements of Web transfers [6]. Between file transfers, the TCP connections were idle for a “think time” also based on the same data (but with the mean reduced by a factor of 10 to generate a heavier load). The only results reported from simulations with these traffic conditions, however, were for buffer occupancy in the RED router which again demonstrated a tendency to stabilize around the maximum threshold for larger numbers of active flows.

Recent work at INRIA has used analytic models and simulation [20] along with live testing on a commercial RED implementation [19] to quantify the performance effects of RED. The emphasis was on quantifying how RED influences loss rates, patterns of consecutive loss, mean delay, and delay jitter for mixes of “bursty” (TCP) and “smooth” (UDP) traffic, when compared with tail-drop FIFO queue management. The results from analytic models were confirmed with ns simulations for a number (up to 300) of continuously sending TCP connections sharing a bottleneck link with UDP flows operating at 10% of the link capacity. They concluded that TCP “goodput” does not improve significantly with RED and this effect is largely independent of the number of flows. They also observed that the mean queueing delay is lower with RED but has a much larger delay variance. In essence, the RED router behaved as a tail-drop router with a queue length equal to the maximum threshold.

Even though the INRIA work considers the effect of both queueing delay and drop rates at routers, it does not integrate these effects with the dynamics of TCP congestion control and retransmission to determine the overall result on end-to-end response times for interactive or Web-like traffic. Moreover, the goal in these experiments was to explore how changes in Cisco’s WRED configuration parameters could be used to control performance. The measures of performance were throughput, bytes sent, and percentage of UDP drops. There were no measurements of delays or end-to-end response times. Their conclusion was that determining the best combinations of RED parameters is difficult and, overall, RED did not show much better performance than tail-drop FIFO (except with larger queue sizes where RED did show some improvement in performance).

We are aware of only two available reports from network operators that have conducted pilot tests of RED in production:

those by Doran at Ebone [7] and Reynolds at QualNet (now Verio) [24]. Doran’s measurements using the Cisco implementation indicate that RED was able to sustain near 100% utilization on a 1920 kb/s customer-access link where tail-drop FIFO could not. Reynolds used the Cisco implementation of WRED on both a DS3 core network link and a DS1 customer-access link. For the heavily congested periods on the core link, it was found that a wide separation of queue thresholds ( $\min_{th} = 60$ ,  $\max_{th} = 500$ ) produced the best tradeoff for link utilization and low drop rates and was somewhat superior to tail-drop FIFO. The default values for drop probability (1/10) and smoothing factor (1/512) were used and their effects not studied. For the customer access DS1 links, (apparently) the default settings were used. These links were congested only during some intervals and some increase in end-to-end latency was observed with RED but the claim was made that “... *the user is not, in my opinion, inconvenienced, and has the benefit of limited packet loss...*” [24].

In summary, while the results from these studies have added important pieces of evidence to the growing corpus of information about RED, important elements are missing. In particular, none of the work we found explicitly considers RED interactions with Web-like traffic where end-to-end response time is the primary measure of performance. Further, many of the results on RED performance are based on “best case” simulations in which a constant number of TCP connections, each sending continuously, share a queue facing a bottleneck link. In the work reported here, we consider the opposite “worst case” in which there is a dynamically changing number of TCP connections with highly variable lifetimes.

### III. EXPERIMENTAL METHODS

#### A. Experimental Network

For our experiments, we constructed a laboratory network that models an enterprise or campus network having a single wide-area link to an upstream Internet service provider (ISP). All traffic using the ISP link is Web traffic where the requesters (browsers) are all located on the enterprise or campus network and all the requests are satisfied by Web servers located somewhere on the Internet beyond the ISP link.

The laboratory network used to emulate this configuration is shown in Fig. 2. All systems shown in this figure are Intel architecture machines running FreeBSD 2.2.8. At one edge of this network are machines that run instances of a Web request generator (described below), each of which emulates the browsing behavior of hundreds of human users. At the other edge of the network are another set of machines that run instances of a Web response generator (also described below) that creates traffic in response to the browsers’ requests. In the remainder of this paper, we refer to the machines running the Web request generator simply as the “browser machines” (or “browsers”) and the machines running the Web response generator as the “server machines” (or “servers”). The browser and server machines have 10/100 Mb/s Ethernet interfaces configured to run at only 10 Mb/s and are attached to a switched VLAN on a Cisco Systems Catalyst 5000. All browser machines are on one VLAN and all server machines are on a separate VLAN.

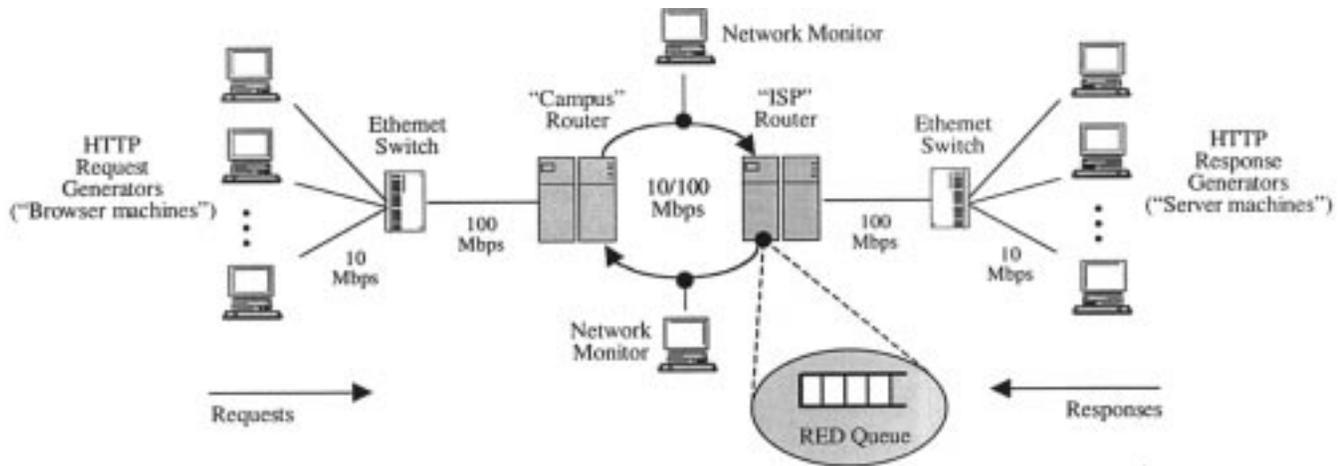


Fig. 2. Experimental laboratory network diagram.

At the core of this network are two router machines running the ALTQ version 1.2 extensions to FreeBSD. ALTQ extends the network-interface output queuing discipline to include FIFO, RED, CBQ, and WFQ queue management [16]. These router machines are 300-MHz Pentium II's. Each router machine has one 100 Mb/s Ethernet interface attached to one of the switched VLAN's on the Catalyst 5000. Each router machine also has two additional 10/100 Mb/s Ethernet interfaces configured to create two point-to-point Ethernet segments (using two hubs) that connect the routers as shown in Fig. 2. Static routes are configured on the routers so that all traffic flowing from the servers to the browsers uses one Ethernet segment and all traffic flowing in the opposite direction uses the other Ethernet segment. This configuration allows us to approximate the full-duplex behavior of the typical wide-area link to an ISP from a customer's network. By configuring the router-to-router Ethernet segments to run at only 10 Mb/s, we can make our representation of the ISP link be a potential bottleneck since the aggregate bandwidth available to the machines at each edge of the network is constrained only by the 100 Mb/s links from the VLANs to the routers. When the links connecting the routers are configured to run at 100 Mb/s, the bottleneck is removed.

Another important factor in modeling this configuration is the effect of end-to-end latency. We use the *dumynet* [8] component of FreeBSD to configure in-bound packet delays on the end systems to emulate different round-trip times between each pairing of a browser machine and a server machine. The delays ranged from 7–137 ms and were derived from measurement data obtained at the *NetStat.net* Web site [21]. The delays were chosen to represent a sample of Internet round-trip times within the continental U.S. A given delay represents the minimum round-trip time experienced by an arbitrary TCP connection between a given pair of client and server machines in our experiments (assuming no delays in the two routers). As explained below, the distribution of TCP connections over pairs of machines should be approximately uniform and, thus, we can calculate the mean minimal round-trip time for all TCP connections sharing the network as approximately 79 ms. The default TCP window size in FreeBSD of 16K bytes was used on all the end systems. (For other characteristics of the TCP implementation as well as the actual delay values used, see [5].)

The instrumentation used to collect network data during runs of the experiments consists of two monitoring programs. One monitor is on the router interface where we are examining the effects of queue algorithms. It calculates a mean and variance of the queue size sampled every 3 ms. The maximum and minimum queue size seen in any sample is also collected. These statistics are logged every 100 ms along with more general information about the number of transmitted and dropped packets. The second monitoring program runs on a separate machine connected to the hubs forming the links between the routers (see Fig. 2). Using a modified version of the *tcpdump* utility, the machine collects the TCP/IP headers in each frame traversing the links and processes these to produce a log of link throughput over each specified time interval (typically one second). End-to-end performance measures such as response times are measured on the end systems as described below.

### B. Web-Like Traffic Generation

The traffic that drives the experiments described here is based on the model of Web browsing developed by Mah [18]. Mah's model is an application-level description of the critical elements that characterize how HTTP 1.0 [22] protocols are used. It is based on empirical data and is intended for use in generating synthetic Web workloads. The data were extracted from more than 230 hours of traces collected on the UC-Berkeley campus in late 1995 and include over 1.6 million HTTP protocol packets. These data were used to compute empirical distributions describing elements necessary to generate synthetic HTTP workloads. The elements of the HTTP model are:

- HTTP request length in bytes;
- HTTP reply length in bytes;
- number of embedded (file) references per page;
- time between retrieval of two successive pages (user "think" time);
- number of consecutive pages requested from a server.

The empirical distributions for all these elements are used in synthetic-traffic generator programs we wrote. The elements that have the most pronounced effects on generated traffic are the size of server responses, the number of requests necessary to download a page (including all embedded references), and the user "think" time between successive page requests. (See

[5] for a more detailed discussion of the Mah data.) We used the Mah model to write Web-traffic generating programs using the normal *socket* system calls provided in FreeBSD. Most of the behavioral elements of Web browsing are emulated in the client-side request-generating program. Its primary parameter is the number of browsing users (typically several hundred) the program is to represent. For each user, the program implements a simple state machine that represents the user's state as either "thinking" or requesting a Web page. If requesting a Web page, a separate TCP connection, as implied by the HTTP 1.0 protocol, is made to the server-side portion of the program for the primary page and each embedded reference (the distribution of embedded references per page is used to generate a random value). Another parameter of the program is the number of concurrent TCP connections allowed on behalf of each browsing user to make embedded requests within a page (this parameter is used to mimic the behavior of Netscape and Internet Explorer).

For each request, a message of random size (sampled from the request size distribution) is sent to the server program. This message contains a value that represents the number of bytes the server is to return as a response (a random sample from the distribution of response sizes). The server sends this number of bytes back to the browser and closes the TCP connection. For the experiments reported here, the server's "service time" is set to zero so the response begins as soon as the request message has been received and parsed (this roughly models the behavior of a Web server or proxy having a large main-memory cache with a hit-ratio near 1.0). For each request-response pair, the browser program logs its response time. Response time is defined as the elapsed time in milliseconds between the time of the socket `connect()` operation and the time the response is completed and the connection is closed. Note that this response time is for each element of a page, not the total time to load all elements of a page.

When all the request-response pairs for a page have been completed, the emulated browsing user enters the "thinking" state and makes no more requests for a period of time sampled from the think-time distribution. The number of page requests the user makes in succession to a given server machine is sampled from the distribution of consecutive page requests. When that number of page requests has been completed, the server to handle subsequent requests is selected randomly and uniformly from the set of active servers. The number of emulated users is constant throughout the execution of each experiment.

The HTTP 1.0 protocol implies the use of a new TCP connection for each request-response pair. This protocol is gradually being replaced by the more efficient HTTP 1.1 protocol which allows multiple and pipelined requests to reuse TCP connections [22]. While some data have been reported (e.g., [10]) suggesting that as many as 30% of HTTP requests now use the HTTP 1.1 protocol, we have been unable to find data or models sufficient for building a synthetic workload generator for HTTP 1.1. For these reasons we generate only HTTP 1.0 traffic in our experiments. We note, however, that the older HTTP 1.0 protocols are expected to represent a very significant portion of Web traffic for some time because of difficulties with migrating the installed base of browsers. Furthermore, our focus on HTTP 1.0 serves as a worst-case analysis of RED performance.

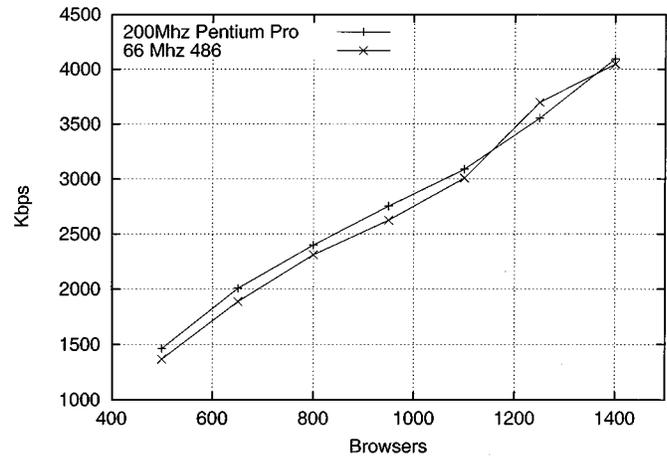


Fig. 3. Offered load as a function of the number of simulated users on one machine.

### C. Experiment Calibrations

There are two critical elements of our experimental procedures that had to be calibrated before performing experiments: 1) ensuring that no element on the end-to-end path represented a primary bottleneck other than when the links connecting the two routers are limited to 10 Mb/s; and 2) the offered load on the network can be predictably controlled using the number of emulated browsing users as a parameter to the traffic generator. To perform these calibrations, we first configured the two segments connecting the routers to eliminate congestion by running at 100 Mb/s.

The first calibration performed was to verify that the traffic generator programs did not have any resource constraints that limited their ability to emulate hundreds of users. These programs were implemented using efficient programming techniques for managing large numbers of socket connections (based in part on Banga and Druschel's scalable methods for generating HTTP requests [3]). For this calibration we first selected the slowest machine in our network (a 66 MHz 486) to run the browser program. We ran one instance of the server-side program on each of the server machines and configured the browser program to select uniformly from all servers for each new sequence of page requests. The number of browsing users was varied from 500 to 1400 and the bandwidth used on the 10 Mb/s interface to the browser machine is plotted in Fig. 3 as a function of the number of simulated browsing users. These results show that over this range of users, there is a linear increase in generated traffic and the traffic is significantly less than the capacity of the host's 10 Mb/s interface. We repeated this experiment with a 200 MHz Pentium Pro with the results also shown in Fig. 3 for further confirmation that CPU and interface speeds of the end system are not resource constraints. Thus if traffic generation machines are limited to simulating no more than 1400 users each, we can be confident that the number of users simulated in an experiment is accurate and reproducible.

A second concern is that a single program can not faithfully simulate hundreds of browsers because by default, a single FreeBSD process can use at most 64 sockets simultaneously. However, because user think times are much longer than the times required to request pages, most of the emulated users

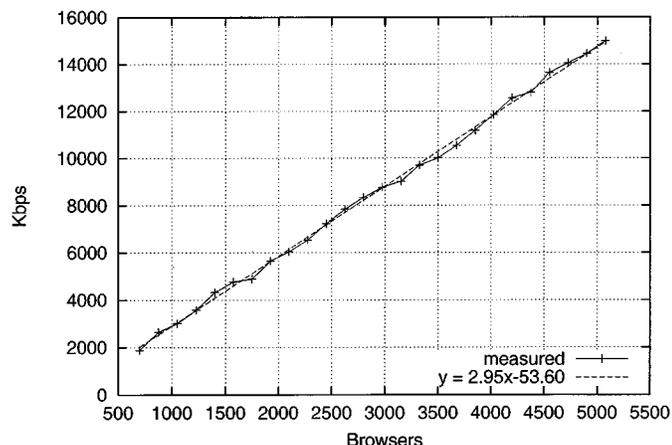


Fig. 4. Offered load as a function of the number of simulated users on seven machines.

are idle at any time. We explicitly performed experiments to demonstrate that the 64 socket descriptors limitation was never encountered in practice. With a similar experiment we also verified that even the slowest server machine could handle a maximum number of expected requests without reaching a resource limitation.

For the next calibration, we ran an instance of the browser program on each of the browser machines and again uniformly distributed requests across all server machines. Each browser was configured to emulate the same number of users with the total users varied from 700 to 5075. Fig. 4 plots aggregate traffic on the path carrying response traffic from the servers as a function of emulated browsers (users). Again the load is a linear function of browsers indicating there are no fundamental resource limitations in the system and generated loads can easily exceed the capacity of a 10 Mb/s link.

With these data, we can determine the number of emulated browsers that would generate a specific offered load in this configuration if there were no bottleneck link present. This capability is used in subsequent experiments to control the offered loads on the network, including loads that nominally exceed the capacity of a 10 Mb/s link. For example, if we want to generate an offered load equal to the capacity of a 10 Mb/s link, we use Fig. 4 to determine that we need to emulate approximately 3400 browsing users; for a load of 110% (11 Mb/s) we need to emulate 3750 users.

A motivation for choosing Web-like traffic to drive these experiments was the assumption that properly generated traffic would exhibit highly variable and bursty demands on the network. To illustrate that this is indeed realized with our experimental setup, we have plotted the results from one of the calibration experiments (3500 browsers) in Figs. 5 and 6. These plots show the number of requests initiated during each one second interval (each request requires a new TCP connection) and the number of bytes requested (not necessarily received) in each one second interval. Clearly, these show the highly bursty nature of the traffic actually generated.

#### D. Experimental Procedures

Each experiment was run using the following procedure. After initializing and configuring all router and end-system

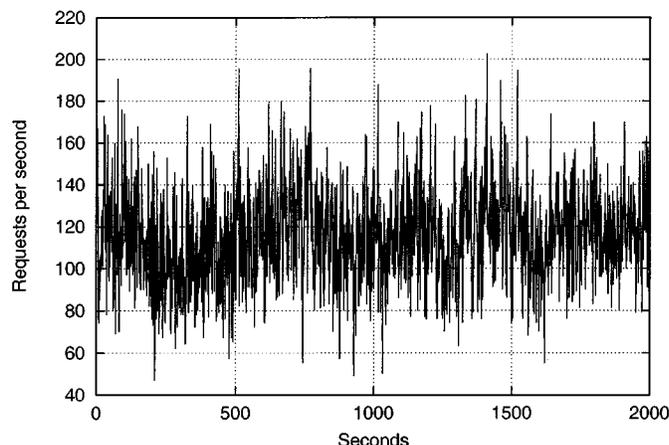


Fig. 5. Requests per second from 3500 users.

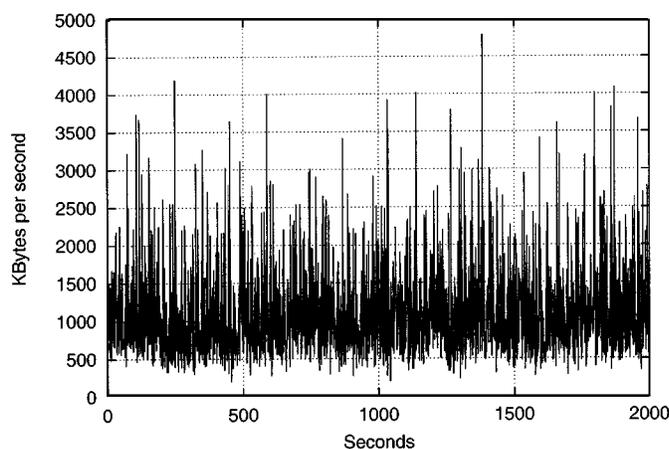


Fig. 6. Bytes requested per second from 3500 users.

parameters, the server-side processes were started followed by the browser processes. Each browser emulated an equal number of users chosen, as described above, to place a nominal offered load on an unconstrained (100 Mb/s) network. The offered loads used in the experiments were chosen to represent 50, 70, 80, 90, 98, or 110 percent of the capacity of a 10 Mb/s link connecting the two router machines. Loads exceeding 110% were tried; it turned out, however, that the extreme duration of the connections when using a congested link caused the traffic generators to occasionally use all available sockets and fail to generate the desired level of traffic. Because the measured response times at a load of 120% had deteriorated well beyond levels that most users would tolerate, we decided to not consider loads beyond 110% on the congested link.

Each experiment was run for 90 minutes, but data collected during the first 20 minutes was discarded to eliminate startup and stabilization effects. These effects are illustrated in Fig. 7 which shows a plot of mean response times for requests during each one second interval in a typical experiment. Fig. 8 gives a plot of the cumulative distribution of response times at a load from 3500 browsers in an unconstrained network. Note that about 90% of the requests complete in 500 ms or less. Fig. 8 represents the best-case performance for HTTP request-response pairs and will be used as a basis for comparison with experiments on the constrained (congested) network link. Table II

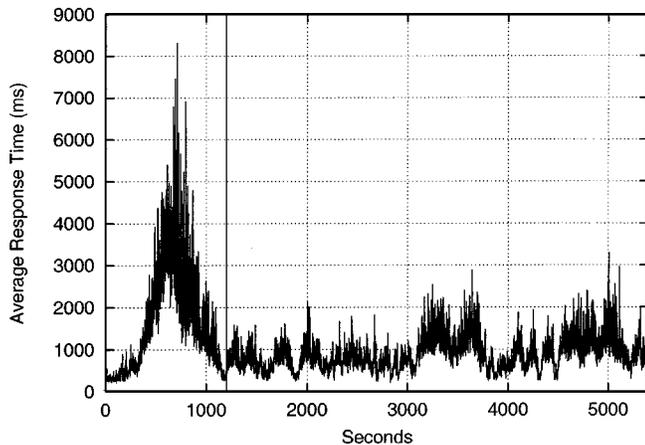


Fig. 7. Average response time per second during an experiment. The plot includes the initial 20 minutes, where the traffic generators are started and stabilize.

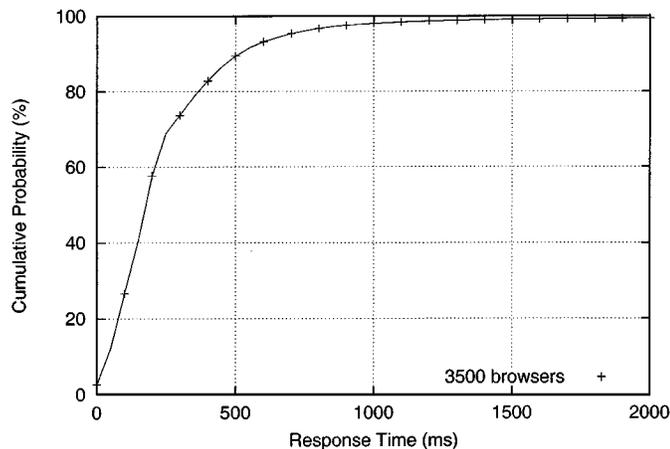


Fig. 8. Cumulative response time distribution for 3500 users on the unconstrained (100 Mb/s) network.

TABLE II  
TYPICAL NUMBERS OF REQUESTS IN A 70-MIN INTERVAL

Load %	Requests	Load %	Requests
50	240,379	90	425,293
70	329,638	98	461,837
80	375,673	110	521,561

shows the number of requests generated during a 70-min interval for each of the loads in typical runs on the unconstrained network.

Because responses are much larger than requests, the load on the link between routers that carries traffic from the servers to the browsers (the link from the “ISP” router to the “campus” router in Fig. 2) will be much greater than that on the link carrying traffic in the opposite direction. Consequently, only the effects of different queue management algorithms on the IP output queue for this link interface are reported here. The IP output queues for the link interfaces on all other machines in the network were tail-drop FIFO queues with the FreeBSD default queue size of 50 elements. Data collected on these interfaces using the *netstat* function showed no dropped packets.

The key indicators of performance we use in reporting our results are the end-to-end response times for each request–response pair. We report several measures of response times including the median, the percent of requests completing in intervals of 0–1, 1–2, 2–3, and greater than three seconds, and plots of the cumulative distributions of response times (usually showing only times less than or equal to two seconds). We also measured the percent of IP datagrams dropped at the link queue, the mean queue size, and the link throughput actually achieved on the bottleneck link.

#### IV. FIFO RESULTS

To establish a baseline for evaluating the effects of using RED on interface queues for links carrying only Web traffic, we first examined the effects of FIFO queues with tail-drop behavior in our experimental network. For these experiments we created a bottleneck between the two routers by configuring the two segments connecting the router machines to run at 10 Mb/s using 10 Mb/s hubs (see Fig. 2). The critical parameter for a FIFO queue is the size of the buffer space allocated to hold the queue elements. Guidelines (or “rules of thumb”) for determining the “best” queue size have been widely debated in various venues including the IRTF *end2end-interest* mailing list [9]. The guideline that appears to have attracted a rough consensus is to provide buffering approximately equal to 2–4 times the bandwidth-delay product of the link. Bandwidth in this expression is that of the link for the interface using the queue and delay is the mean round-trip time for all connections sharing the link—a value that is, in general, very difficult to determine. For our experimental network, the mean minimum RTT can be computed as 79 ms and the 10 Mb/s link has a bandwidth-delay product of approximately 96K bytes [5]. FreeBSD queues are allocated in terms of a number of buffer elements (*mbufs*) each with capacity to hold an IP datagram of Ethernet MTU size. We measured the mean IP datagram size in our generated Web response traffic to be just over 1K bytes so the FIFO queue should have approximately 190–380 queue elements to fall within the guidelines.

We ran a number of experiments with a FIFO queue on the bottleneck link varying the offered load and queue size. Fig. 9 shows the cumulative response time distributions for different FIFO queue sizes at loads of 80%, 90%, 98%, and 110%. At a load of 80%, there is little effect from increasing the queue size from 30 to 240 elements. At 90% load, we begin to see queue size having more significant effects on response times and observe that a queue size of 120 elements is a reasonable choice for this loading. The effect that queue size has on response times depends on the size of the HTTP response data as is shown in the plots for 98% load. Increasing the queue size from 30 to 120 has a slightly negative effect on relatively short responses that could complete in a few hundred milliseconds by increasing the amount of time each packet spends in the queue. For a 10 Mb/s Ethernet link and an average frame size around 1 KB, approximately 1000 packets can be forwarded per second. Thus a packet arriving at the queue already containing 100 packets has to wait approximately 100 ms on the router. Such a delay is significant for requests with short responses that may otherwise complete

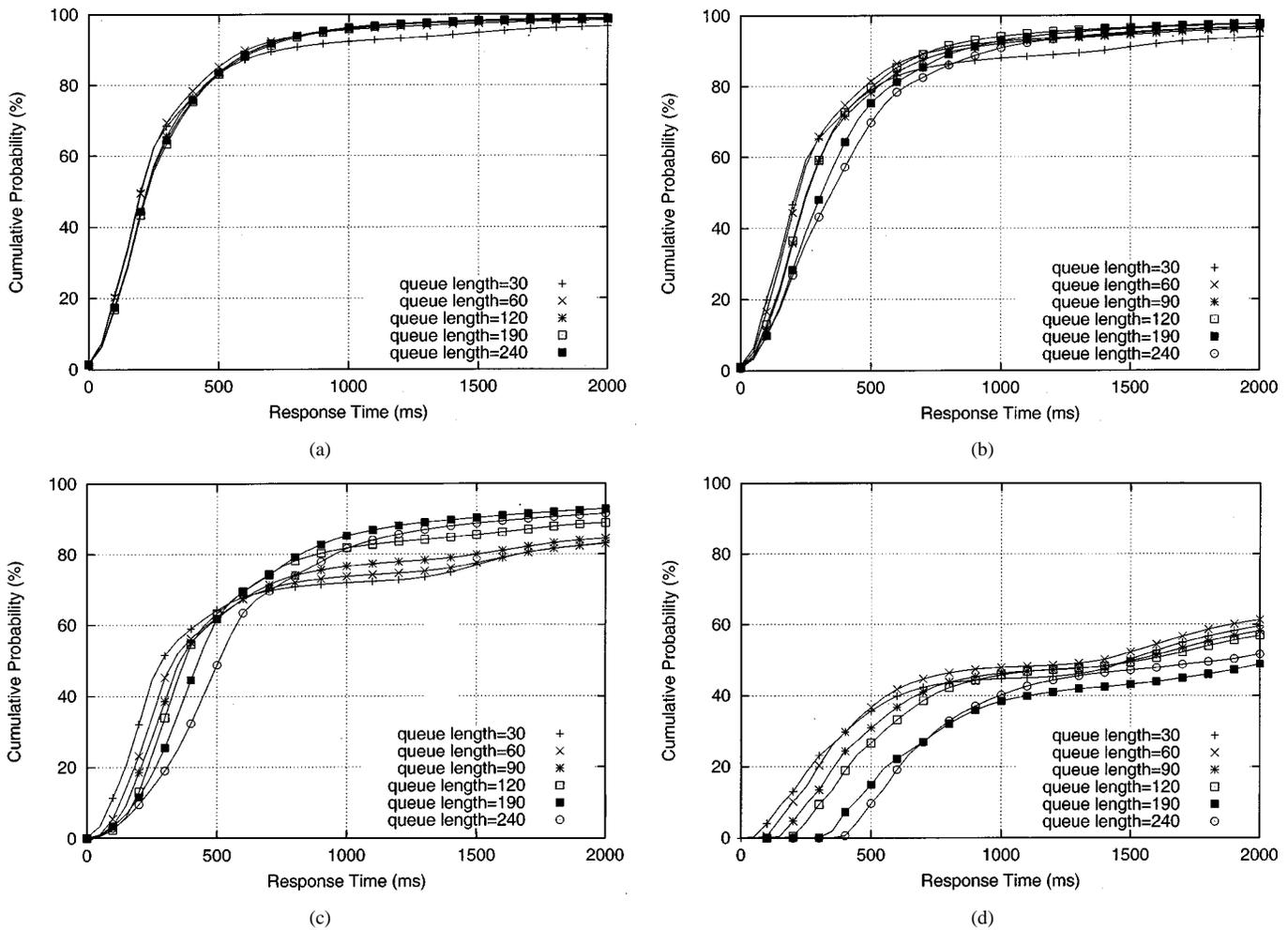


Fig. 9. (a) FIFO performance at 80% load. (b) FIFO performance at 90% load. (c) FIFO performance at 98% load. (d) FIFO performance at 110% load.

within 200–350 ms. On the other hand, increasing the queue size from 30 to 120 reduces response times significantly for long requests. Even though the time spent in the queue by each packet is longer, the reduced rate of drops means that longer responses are less likely to encounter retransmission timeouts (which are often longer than queueing delays by a factor of 5–10 times). At queue sizes of 190 or 240, the increase in response times for short requests appears to offset any improvement gained for longer requests from reduced drops.

Our results indicate that, overall, a FIFO queue size of 120 elements (about 1.25 times the bandwidth-delay product) to 190 elements (2 times bandwidth-delay) is a reasonable choice for loads up to the link capacity. For offered loads that only slightly exceed the link capacity (e.g., 110%), we observe that queue sizes beyond 120 only exacerbate an already bad situation. Additional measures of performance in these experiments, including link utilization and drop rates, confirm that our selection of queue sizes of 120–190 represent reasonable tradeoffs for response times without significant loss of link utilization or high drop rates [5].

These experiments illustrate (as queuing theory predicts) the dramatic effect that offered loads near or slightly beyond the link capacity have on response times. Fig. 10 shows the cumulative distribution of response times for these loads with a FIFO

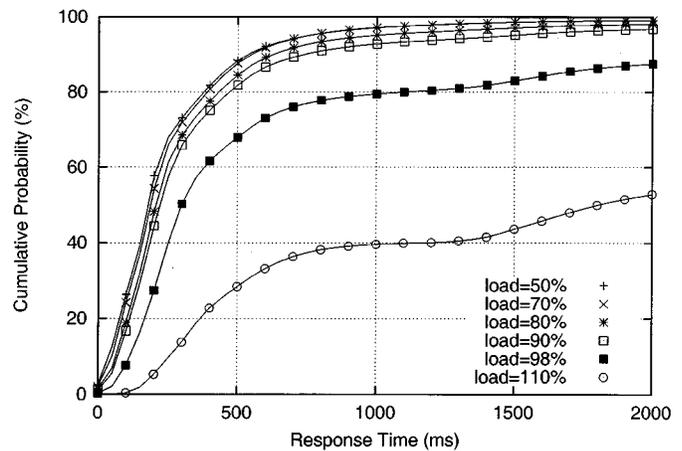


Fig. 10. FIFO performance for different loads with a queue length of 120 elements.

queue of 120 elements. Clearly, response times degrade sharply when the offered load approaches or exceeds link capacity. If an ISP has links that experience utilization above 90% over intervals greater than a few minutes, response time for Web users are seriously impacted. A second important observation is that at loads below 80%, there is no significant change in response times as a function of load.

## V. RED RESULTS

The goal for our experiments with RED was to determine parameter settings that provide good performance for Web traffic. We also examined the tradeoffs among the different parameters in tuning for performance. The RED queueing mechanism has five different parameters for adjusting the algorithm's behavior. An exhaustive search for the best parameter values is impossible because of the number of possible combinations of values. Our approach for the RED experiments was to design an initial set of experiments that could give a broad approximation of parameter values that result in good HTTP performance. We then examine the effects of varying each parameter individually using this initial determination as a baseline.

From our experiments with FIFO, it is clear that there is a complex tradeoff between response times for short responses that can be completed in a few hundred milliseconds (best with a short queue) and response times for longer responses (best with longer queues and lower drop rates). The original Floyd and Jacobson paper [13] suggests guidelines for tuning parameters that have been revised based on subsequent experience and analysis (see [15] for the current guidelines). These guidelines suggest that the most fundamental effects are determined by the  $\min_{th}$  and  $w_q$  parameters which control tradeoffs between average queue size and sensitivity to the duration of periods of congestion. For our initial experiments we decided to eliminate the size of the physical queue as a factor and set the number of queue elements to 480, more than double the largest average queue size seen in the FIFO experiments. In these experiments we varied  $\min_{th}$  beginning with the guideline value of 5 and ranging up to 120. We fixed  $\max_p$  at 0.10,  $w_q$  at 0.002 (actually 1/512), and  $\max_{th}$  at 3 times  $\min_{th}$  as suggested in the current guidelines.

Each of the parameter settings was tried at six different offered loads: 50%, 70%, 80%, 90%, 98%, and 110%. At 50% load the number of dropped packets was between 0.00% and 0.01% of the total number of packets transmitted. This means that at loads of 50% and below, there is limited room for increasing the performance of the router queueing mechanism. Post processing of the logs shows that the queue size never reaches the maximum value of 480 even at a load of 110%, though it is possible in a worst-case scenario. As expected, the performance changes significantly as the load is increased from 50% to 110%. Fig. 11 illustrates typical results from these experiments by showing the effect of varying loads on response time distributions with  $(\min_{th}, \max_{th})$  set to (30, 90).

It is encouraging to see that performance degradation only occurs at loads greater than 70%, especially when combined with the fact that the drop rates at 50% load never exceeds 0.01% of the packets received at the router. This indicates that parameter tuning will have limited effect until loads reach levels of 70%–80% of link capacity. When loads exceed 70%, the performance decreases monotonically as the load increases. The most significant performance decrease occurs at load levels of 90%–110%. These are the most interesting targets for optimization, since this is where there is significant performance to gain.

We start by exploring possible choices for  $\min_{th}$  and  $\max_{th}$ . Fig. 12 shows the response time distributions for the 90% and

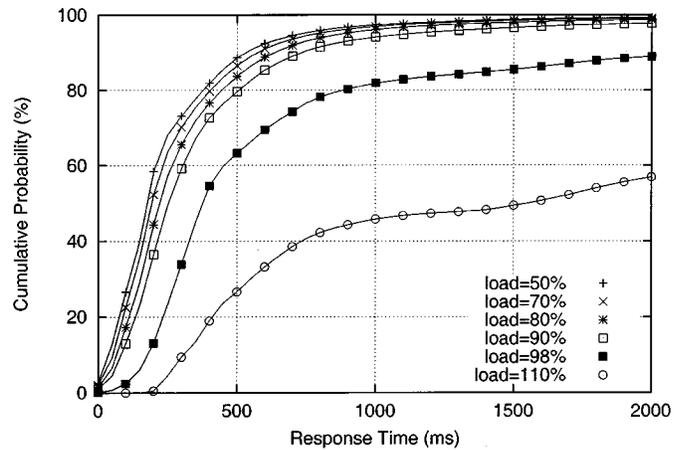


Fig. 11. Performance of RED at different loads.  $w_q = 1/512$ ,  $\max_p = 1/10$ ,  $\min_{th} = 30$ ,  $\max_{th} = 90$ ,  $qlen = 480$ .

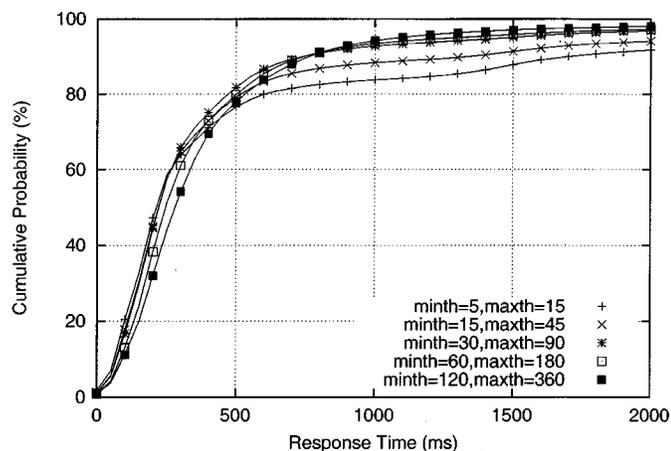
98% offered loads, respectively. These results clearly show that a naive application of the guidelines in [15] with a  $\min_{th}$  of 5 would result in poor performance for Web-dominated traffic. The best overall response-time performance is obtained with values for  $(\min_{th}, \max_{th})$  of (30, 90) or (60, 180). We see, as in the case of FIFO, that there is a tradeoff between better response times for short responses at (30, 90) and improving response times for longer ones at (60, 180), especially at the 98% load. Although the differences are not great, we prefer (30, 90) on the grounds that about 70% of the requests experience somewhat better response times than with (60, 180). One could also argue that (60, 180) is best because it improves the most noticeable delays. (The complete results, including link utilization and drop rates for loads of 90%, 98%, and 110% can be found in [5]. These data indicate a slight drop in link utilization for the (30, 90) setting over the (60, 180) setting.) Like the FIFO results, response times at loads of 110% are quite bad and are not improved by changing the RED settings for  $(\min_{th}, \max_{th})$ .

We next consider varying the ratio between  $\min_{th}$  and  $\max_{th}$  by holding one constant and varying the other. To see the effect of  $\min_{th}$ , we first fixed  $\max_{th}$  at 90 and varied  $\min_{th}$ . We then held  $\min_{th}$  constant at 30 and varied  $\max_{th}$ . We fixed  $\max_p$  at 0.10,  $w_q$  at 0.002 (actually 1/512), and  $qlen$  at 480 as in the previous experiments. Fig. 13 illustrates the effect from varying  $\min_{th}$  on the response time distributions for the 90% load. The results obtained by varying  $\max_{th}$  are similar. The results from these experiments, in general, show only marginal changes in response times (or link utilization) and confirmed the notion that the best balance of response times for all sizes of responses with the loads considered here are achieved with  $(\min_{th}, \max_{th}) = (30, 90)$ . (The complete results from these experiments can again be found in [5].)

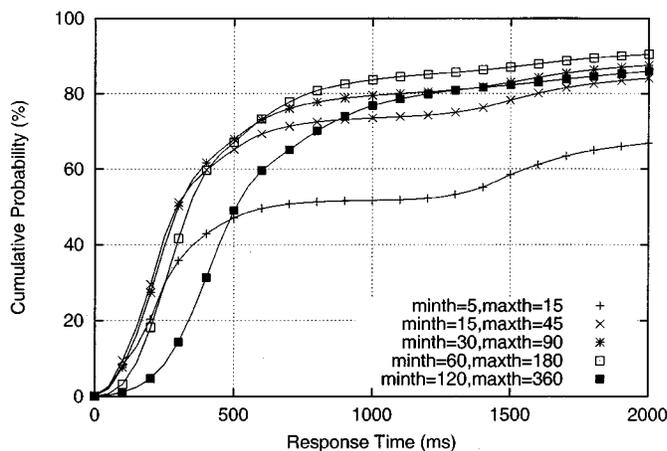
Experiments testing the impact of changing  $w_q$  and  $\max_p$  were combined because of the close relationship between the two parameters. The values used for  $w_q$  were: 1/512, 1/256, and 1/128. (The implementation of RED requires the denominator to be a power of 2.) Decreasing  $w_q$  to 1/1024 was tried, but we found it to be an unrealistic setting that causes reaction to congestion to be quite slow. The values of  $\max_p$  used were 0.05, 0.10, and 0.25. The remaining parameters were fixed at  $\min_{th} =$

TABLE III  
RED PERFORMANCE WITH RECOMMENDED PARAMETERS AND QUEUE LENGTHS

Load %	Queue Length	KB/s	% drop	Mean queue	Median resp.(ms)	% ≤ 1 sec	1<% ≤ 2 sec	2<% ≤ 3 sec	% > 3 sec
90	480	1079	0.8	20.2	266	92.5	4.3	2.0	1.3
90	160	1093	1.1	22.2	278	91.2	4.7	2.4	1.7
90	120	1066	0.7	18.8	266	93.0	4.1	1.7	1.2
98	480	1164	4.1	39.4	345	79.2	8.2	6.3	6.3
98	160	1175	5.9	46.3	397	72.4	9.7	8.2	9.7
98	120	1171	5.5	44.3	377	74.2	9.2	7.7	8.9
110	480	1187	19.7	76.0	1846	39.4	12.9	12.1	35.5
110	160	1188	19.5	76.6	1864	39.1	13.0	12.2	35.7
110	120	1188	18.9	77.0	1840	39.3	13.2	12.5	34.9



(a)



(b)

Fig. 12. (a) Response time CDF for offered load at 90% of link capacity ( $w_q = 1/512, \max_p = 1/10, qlen = 480$ ). (b) Response time CDF for offered load at 98% of link capacity ( $w_q = 1/512, \max_p = 1/10, qlen = 480$ ).

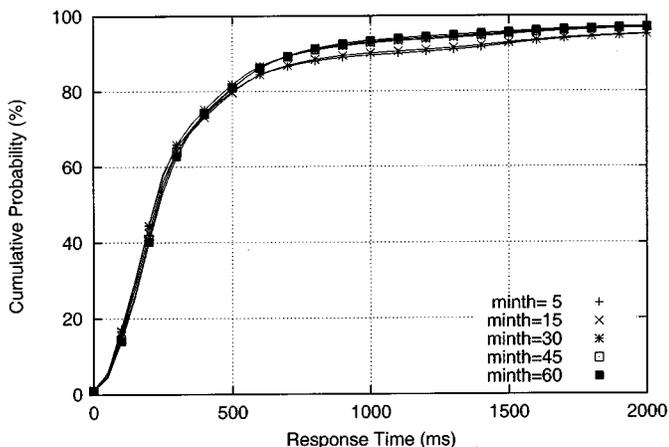


Fig. 13. Effect of changing  $\min_{th}$ . Load = 90% and  $\max_{th} = 90, w_q = 1/512, \max_p = 1/10, qlen = 480$ .

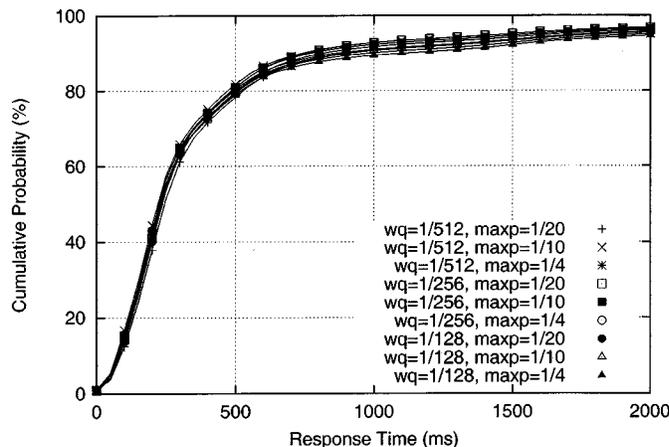


Fig. 14. Results for different values of  $w_q$  and  $\max_p$ . Load = 90%, and  $qlen = 480, \min_{th} = 30, \max_{th} = 90$ .

30,  $\max_{th} = 90$ , and  $qlen = 480$ . All the different settings were tested at loads of 90%, 98%, and 110%.

These experiments showed that at all load levels the setting of  $\max_p$  to 0.25 has a negative impact on performance, because too many packets are dropped. Fig. 14 shows the results from the experiments at 90% load (the results at 98% are similar). At 90% and at 98% load, the difference between the settings occurs beyond the knee (above the 75th percentile) of the CDF, meaning that changes of  $w_q$  and  $\max_p$  mainly impact the longer

flows. Overall, however, we conclude that there is no strong evidence to indicate using values other than the suggested  $w_q = 1/512$  and  $\max_p = 0.10$ .

Finally, we consider the effect of having a limit on the queue size such that there are occasionally forced drops because the instantaneous queue exceeds the buffer space. Table III gives experimental results with our recommended values of RED parameters for actual queue sizes of 480, 160, and 120 elements. These results are very similar to the FIFO results—the 120 ele-

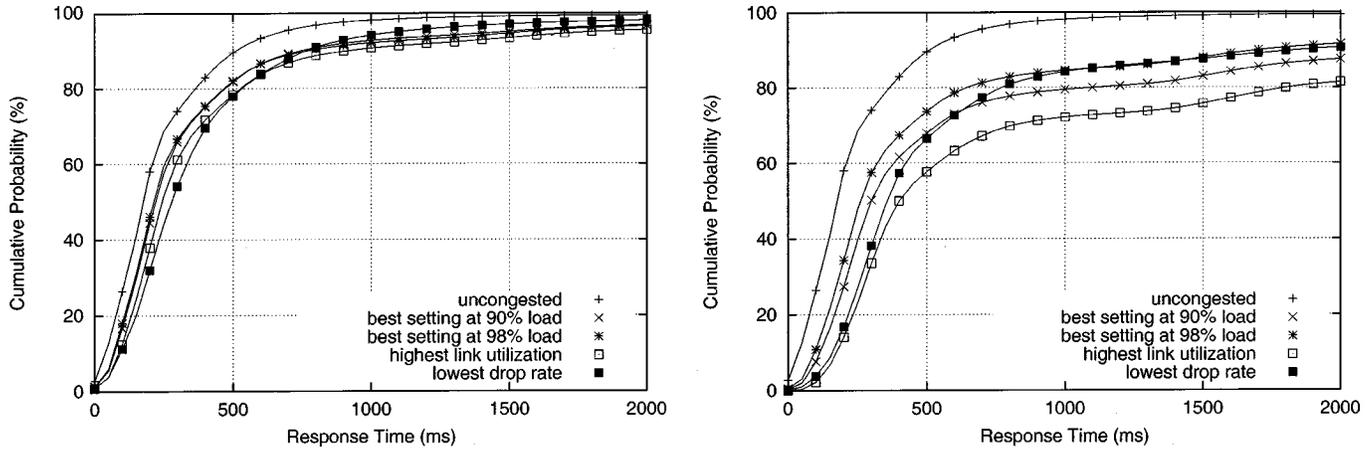


Fig. 15. (a) “Good” RED parameter settings at 90% load. (b) “Good” RED parameters settings at 98% load.

TABLE IV  
EMPIRICALLY DETERMINED “BEST” RED PARAMETER VALUES

Load	$\min_{th}, \max_{th}$	$w_q$	$\max_p$	Notes
90	30,90	1/512	1/10	best overall response
90	30,90	1/512	1/20	highest link utilization
90	120,360	1/512	1/10	lowest drop rate
98	5,90	1/128	1/20	best overall response
98	30,180	1/512	1/10	highest link utilization
98	90,150	1/512	1/10	lowest drop rate

ment queue (1.25 times bandwidth-delay) is a reasonable choice at 90% and 110% loads while a longer queue of 2–3 times bandwidth-delay might provide some advantage at loads just below link saturation.

Our conclusion is that, except for  $\min_{th}$  which should be set to larger values to accommodate the highly bursty character of Web traffic, the guidelines for RED parameter settings and for configuring interface buffer sizes (for both FIFO and RED) also hold for the Web-like traffic used in our experiments. We also conclude that attempting to tune RED parameters outside these guidelines is unlikely to yield significant benefits.

To illustrate the latter point, we examined the entire suite of experiments conducted for the 90% and 98% loads (including some trial experiments with parameter values outside the ranges reported above) to find the combination of settings that gave the best results on three performance measures: “best” response times (a subjective choice because of the tradeoff between improving response times for short versus long responses), best link utilization, and lowest drop rate. These settings are shown in Table IV and the response times shown in Fig. 15. For 90% load, there are relatively small differences between tuning for highest link utilization or lowest drop rates and tuning for response times. At 98% loads, tuning for highest link utilization has potentially serious effects on increasing response times. Note that the “best” overall response times are obtained for the 98% load (only) with parameters that are quite different from our generally recommended settings. (In Fig. 15, the “uncongested” plots refer to the performance on the unconstrained 100-Mb/s network.)

Moreover, there is a significant downside potential for choosing “bad” parameter settings, especially at near-saturation loads. We again searched the entire set of experiments for the 90% and 98% loads looking for combinations of RED parameters that produced response times that (subjectively) represented poor choices (i.e., choices that increased response times significantly for larger numbers of either short or long responses). Fig. 16 shows these results. Clearly, some parameter settings produce results that are considerably less desirable than our recommended ones.

## VI. ANALYSIS OF RED RESPONSE TIMES

While a detailed analysis of the causes of the response time distributions observed under RED remains the subject of further study, we are able to report the results of a preliminary analysis. We repeated two of the RED experiments reported in Fig. 12(b) using more elaborate instrumentation to provide additional data for each HTTP request–response pair. In addition to end-to-end response times, this instrumentation allowed us to determine a detailed breakdown of the number and types of retransmission events occurring in the TCP connections. We repeated RED experiments for an offered load of 98% with two sets of RED parameter values that produced clearly different response time results;  $(\min_{th}, \max_{th}) = (5, 15)$  and  $(60, 180)$ . In both cases, the remaining RED parameters were fixed at  $\max_p = 1/10$ ,  $w_q = 1/512$ , and  $qlen = 480$ .

Table V gives the percentage of TCP connections that experienced no retransmissions, those that experienced one or more retransmissions of a SYN segment, one or more retransmissions of a FIN segment, one or more retransmissions of a data segment, and those with any combination of SYN, FIN, and data segment retransmissions. In this analysis, we consider only retransmissions on the congested path from the server to the client (there were no observed packet drops on the reverse path).

Figs. 17 and 18 give the cumulative distributions of response times for those connections that experienced no retransmissions and for the connections that experienced retransmission events of the types described above. Also shown is the cumulative distribution of response times for all connections. In Fig. 17, we observe that the response times for about 50% of

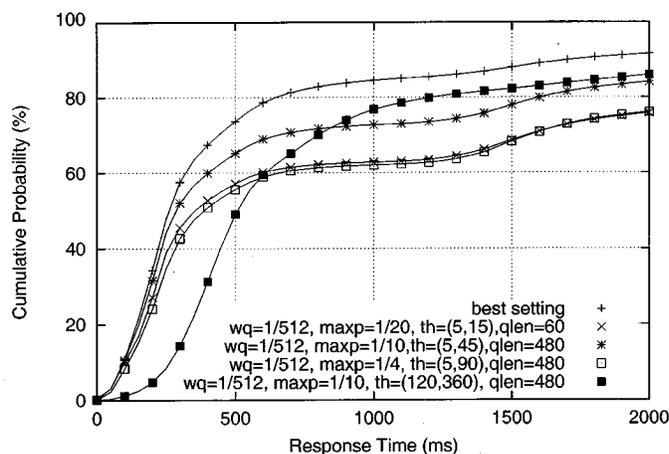
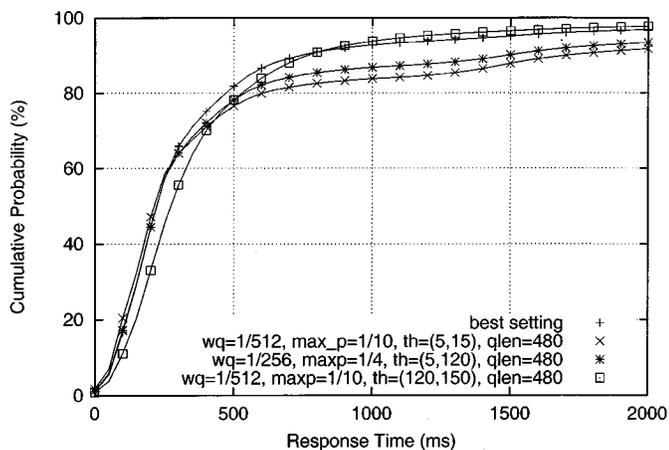


Fig. 16. (a) “Bad” RED parameters settings at 90% load. (b) “Bad” RED parameters settings at 98% load.

TABLE V  
SUMMARY RETRANSMISSION STATISTICS FOR EXPERIMENTS WITH MORE DETAILED INSTRUMENTATION

Class of retransmission event	% of all TCP connections	
	$(min_{th}, max_{th}) = (5, 15) (60, 180)$	
No retransmissions	56.1	87.1
1 or more retransmissions	43.9	12.9
1 or more SYN segments	7.4	2.0
1 or more FIN segments	6.0	2.0
1 or more data segments	25.5	8.5
Combined SYN/FIN/data	5.0	0.4
Total TCP connections	439,979	460,022
Total segments lost	12.4%	2.4%

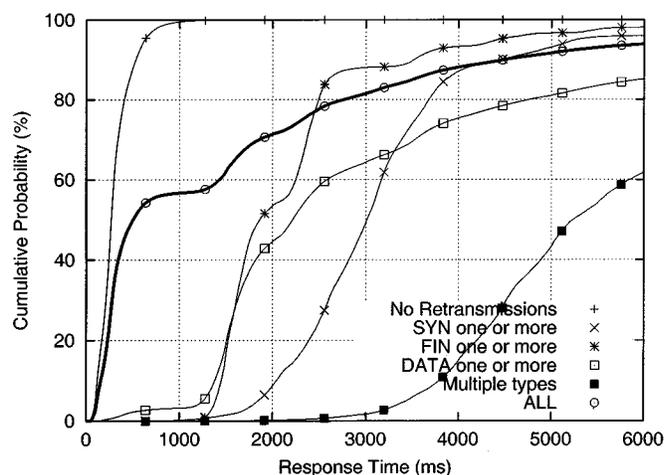


Fig. 17. Absolute performance of flows experiencing retransmissions  $[(min_{th}, max_{th}) = (5, 15)]$ .

the connections with FIN or data retransmissions are shifted relative to those with no retransmissions by an amount corresponding to typical retransmission timeouts in our experiments (approximately 1.5 seconds). The response times for connections with SYN retransmissions are shifted even more because of the longer timeouts on TCP connection establishment. Connections with one or more data retransmissions or with combinations of retransmission types have heavier distribution tails (longer response times) because of the cumulative effects of multiple retransmissions.

Comparing the two figures, we observe that response times for those connections having retransmissions are longer in Fig. 18 by a factor somewhat greater than the additional mean queuing delay for this case (about 55 ms). Our preliminary analysis indicates that changes in response times because of retransmissions are a complex combination of factors that influence the retransmission delays. These include the mean queuing delay (which influences the estimated RTT), the deviations in RTT caused by increased variance in queuing delays (which are magnified by a factor of 4 in the TCP algorithm for computing the timeout), the timer granularity (500 ms), and the minimum timeout value (one second). (For a more comprehensive analysis of these and other factors affecting TCP retransmissions, see [1].)

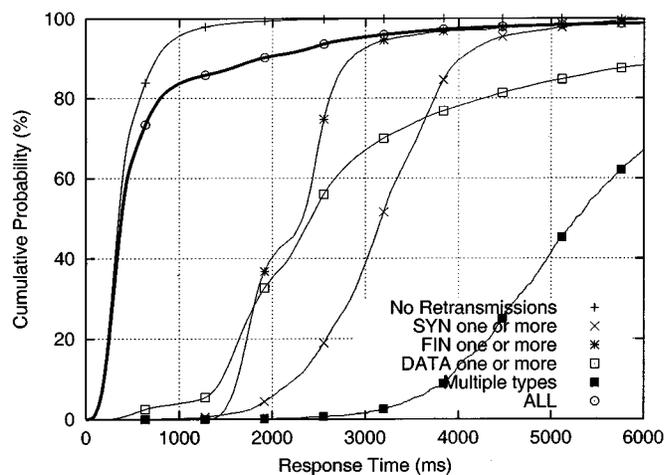


Fig. 18. Absolute performance of flows experiencing retransmissions  $[(min_{th}, max_{th}) = (60, 180)]$ .

The relative contributions of each class of retransmission to the overall response time distribution is shown in Figs. 19 and 20. They show the cumulative distributions of connections with the following characteristics: those with no retransmissions,

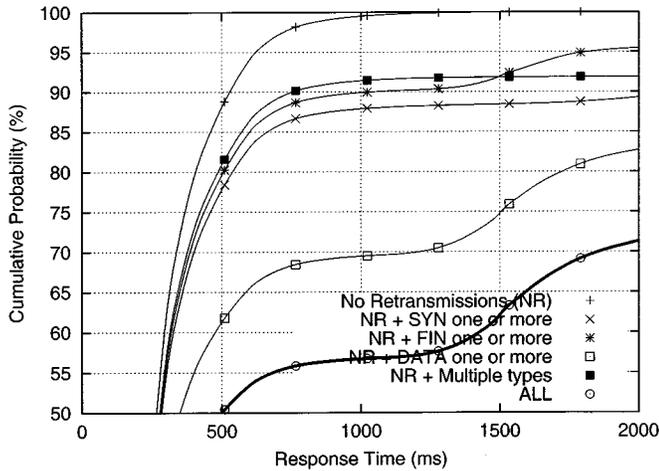


Fig. 19. Relative contribution of flows experiencing retransmissions to total distribution  $[(\min_{t_{th}}, \max_{t_{th}}) = (5, 15)]$ .

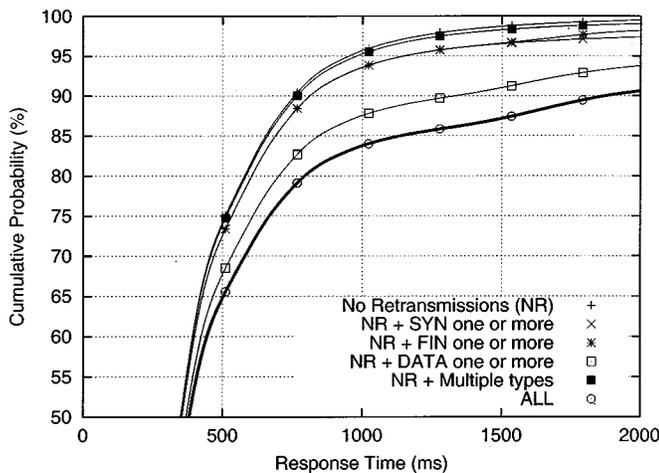


Fig. 20. Relative contribution of flows experiencing retransmissions to total distribution  $[(\min_{t_{th}}, \max_{t_{th}}) = (60, 180)]$ .

those with either no retransmissions or only SYN retransmissions, those with either no retransmissions or only FIN retransmissions, and those with either no retransmissions or only data segment retransmissions. To magnify the relative contributions of each class we only show the portion of the distribution beyond the 50th percentile. Contrary to our expectations, retransmissions of lost SYNs (even when most of the TCP connections transfer relatively few bytes) is far from being the dominant factor leading to increased response times. It is, in fact, data segment retransmissions that have the greatest cumulative effect.

Another view of these dynamics is shown in Fig. 21 that gives a scatter plot of response times versus server reply sizes. There is one dot in this plot for each of the approximately 400 000 connections with reply sizes less than 16 000 bytes in the experiment illustrated in Figs. 17 and 19  $[(\min_{t_{th}}, \max_{t_{th}}) = (5, 15)]$ . Connections experiencing one or more retransmissions are marked with dark black dots while those with no retransmissions are marked with gray dots. Several features of this plot are striking.

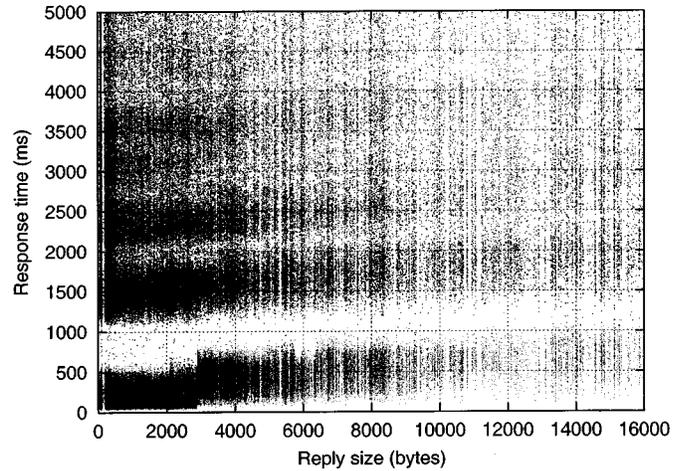


Fig. 21. Scatter plot of response times versus reply size under RED for  $(\min_{t_{th}}, \max_{t_{th}}) = (5, 15)$ .

- the large influence of retransmissions on response times for short responses (e.g., the large number of replies of size less than 4K bytes that take five seconds to complete);
- the clear regions of response times divided between connections with and without retransmissions;
- the distinct bands of response times at intervals roughly proportional to the granularity of the TCP retransmission timer;
- the sharp step increase in response times with no retransmissions for those responses with lengths greater than 2880 bytes (corresponding to the initial TCP congestion window);
- the relatively few connections with retransmissions that avoid a timeout, e.g., with fast retransmission (indicated by black dots in the region dominated by connections with no retransmissions).

This brief analysis has reinforced our view that understanding the effects of RED and FIFO queue management on end-to-end response times for HTTP traffic is a complex issue. It involves many tradeoffs and parameters including not only parameters set on routers but also those controlled at the end systems (e.g., TCP retransmission parameters).

## VII. COMPARING FIFO AND RED

Fig. 22 shows the response time distributions for RED and FIFO with the parameters selected as a result of our experiments at offered loads of 90%, 98%, and 110%, respectively. Also included for reference are the response time distributions at these loads from the calibrations on the unconstrained network. The only case in which there is a distinct advantage from using RED is at the 98% load where response times for shorter responses (80% of requests) are improved with carefully tuned RED parameters. Note that in Fig. 22(a) we include the response time distribution for both  $(\min_{t_{th}}, \max_{t_{th}}) = (30, 90)$  and  $(60, 180)$ . Recall that these settings gave nearly identical performance at 90% load and we selected  $(\min_{t_{th}}, \max_{t_{th}}) = (30, 90)$  as the best setting at this level for largely subjective reasons.

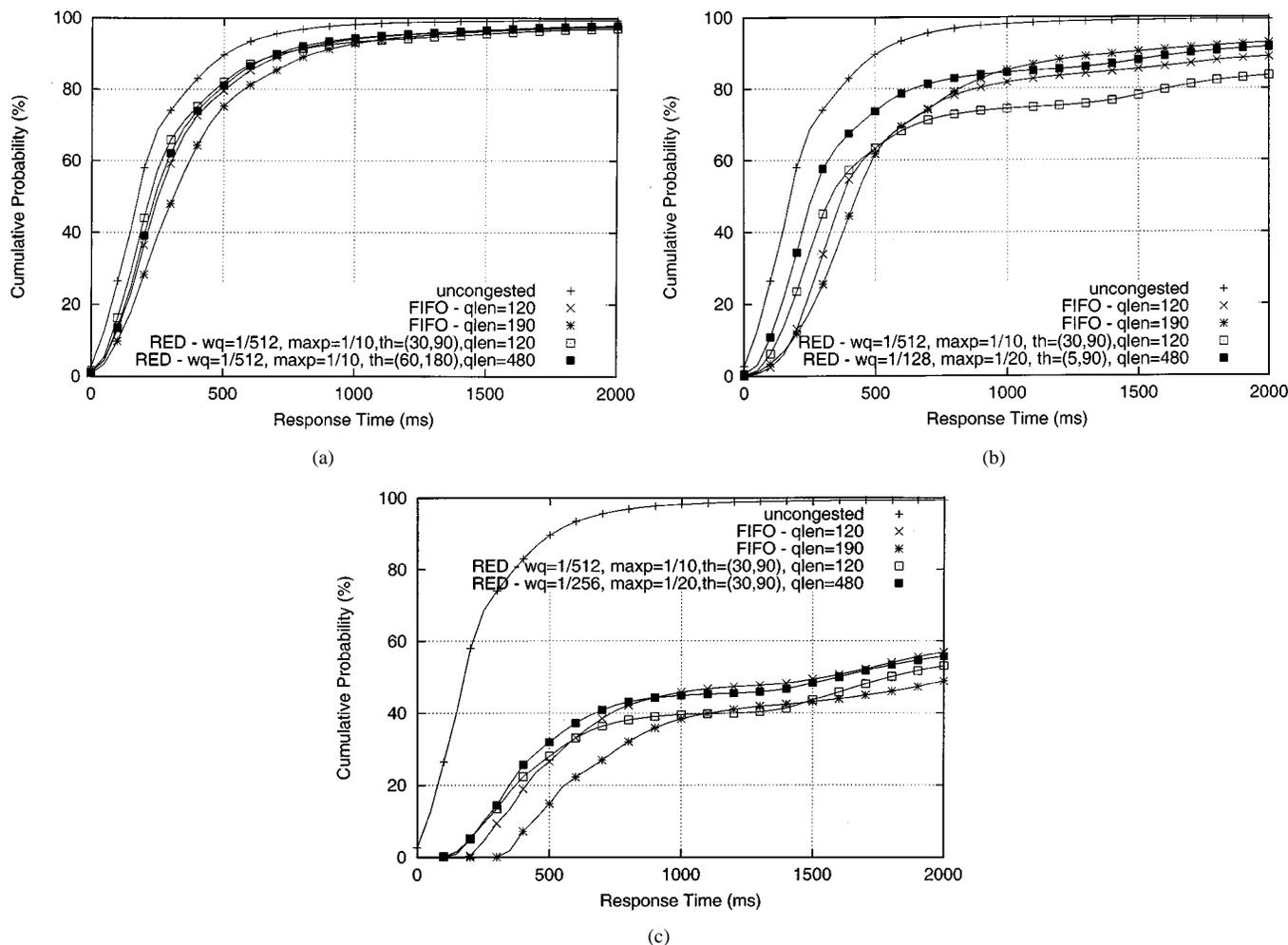


Fig. 22. (a) FIFO and RED at 90% load. (b) FIFO and RED at 98% load. (c) FIFO and RED at 110% load.

### VIII. CONCLUSIONS AND FUTURE DIRECTIONS

Based on our experiments, we summarize our conclusions as follows. Contrary to expectations, for offered loads near or below the levels of link saturation (90% or less), there is little difference in end-to-end response times between the best-tuned RED and tail-drop FIFO configured with 1–2 times the bandwidth-delay product in buffer space. Tuning of the RED parameters generally produces little gain (or loss) in response time performance, however, as illustrated in Fig. 16(a), one can use plausible values for certain RED parameters and produce poorer performance.

At offered loads that approach link saturation (above 90%), RED can be carefully tuned to yield performance somewhat superior to properly configured tail-drop FIFO. The difference is probably significant only between 90% and 100% loading as response times degrade so rapidly above this level that any “improvement” from tuning RED (or FIFO) is, at best, a second-order effect. Moreover, at loads above 90%, response times are more sensitive to the actual values of RED parameters. In particular, there is greater down-side potential from choosing “bad” parameter values as illustrated in Fig. 16(b). This is significant because parameter settings that outperformed

FIFO were arrived at only through extensive trial-and-error experimentation. It was also the case that the RED parameters that provide the best link utilization at this load produce poorer response times.

In general, we observed a complex tradeoff between choosing parameters that improve response time for short responses (those consisting of only a few TCP segments) and those that improve response times for longer responses. We have chosen to favor those parameter settings that improve performance for the largest fraction of responses, and hence have focused on improving response times for the shorter responses.

Qualitatively these conclusions imply that providing adequate link capacity (utilization less than 90%) is far more important for Web response times than tuning queue management parameters. If one decides to deploy RED for any reason, response times for Web-dominated traffic are not likely to be impacted positively and, unless careful experimentation is performed, response times can suffer. Given the current lack of a widely accepted analytic model for RED performance or field-tested engineering guidelines for RED deployment and the complexity of setting RED parameters, there seems to be no advantage to RED deployment on links carrying only Web traffic.

In applying these conclusions, there are some limitations of this study that should be considered.

- We used packet drops as the only “marking” behavior of RED. Explicit marking by RED for ECN-capable TCP implementations is likely to produce better results.
- We examined only HTTP 1.0 protocols. The interaction of RED with a mix of HTTP 1.0 and HTTP 1.1 traffic should also be analyzed.
- We studied a link carrying only Web-like traffic. More realistic mixes of HTTP and other TCP traffic as well as traffic from UDP-based applications need to be examined. Congestion on both paths on a full-duplex link and over multiple router hops should also be considered.

Removing these limitations to produce a broader perspective on RED behavior is the central theme of our ongoing networking experiments.

#### REFERENCES

- [1] M. Allman and V. Paxson, “On estimating end-to-end network path properties,” in *Proc. SIGCOMM*, Sept. 1999, pp. 263–274.
- [2] F. Anjum and L. Tassiulas. Balanced-RED: An algorithm to achieve fairness in the Internet. [Online]. Available: [http://www.isr.umd.edu/TechReports/ISR/1999/TR\\_99-17/TR\\_99-17.phtml](http://www.isr.umd.edu/TechReports/ISR/1999/TR_99-17/TR_99-17.phtml)
- [3] G. Banga and P. Druschel, “Measuring the capacity of a Web server,” in *Proc. USENIX Symp. Internet Technologies and Systems (USITS)*, Dec. 1997, pp. 61–71.
- [4] B. Braden, D. Clark, J. Crowcroft, B. Davie, S. Deering, D. Estrin, S. Floyd, V. Jacobson, G. Minshall, C. Partridge, L. Peterson, K. Ramakrishnan, S. Shenker, J. Wroclawski, and L. Zhang, “Recommendations on queue management and congestion avoidance in the Internet,” RFC 2309, Apr. 1998.
- [5] M. Christiansen, K. Jeffay, D. Ott, and F. D. Smith. Tuning RED for Web traffic (extended version). [Online]. Available: <http://www.cs.unc.edu/Research/dirt>
- [6] M. Crovella and A. Bestavros, “Explaining World Wide Web traffic self-similarity,” Boston Univ. Computer Science Dept., TR-95-015, Oct. 12, 1995. Revised.
- [7] [Online]. Available: <http://adm.ebone.net/~smd/red-1.htm>
- [8] [Online]. Available: [http://www.iet.unipi.it/~luigi/ip\\_dummysnet/](http://www.iet.unipi.it/~luigi/ip_dummysnet/)
- [9] [Online]. Available: <ftp://ftp.isi.edu/end2end/end2end-interest-1998.mail>
- [10] [Online]. Available: [http://www.research.att.com/~anja/w3c\\_webchar/](http://www.research.att.com/~anja/w3c_webchar/)
- [11] W. Feng, D. Kandlur, D. Saha, and K. Shin, “A self-configuring RED gateway,” in *Proc. INFOCOM*, Mar. 1999, pp. 1320–1328.
- [12] ———, “Blue: A new class of active queue management algorithms,” Univ. of Michigan, Ann Arbor, Tech. Rep. CSE-TR-387-99, Apr. 1999.
- [13] S. Floyd and V. Jacobson, “Random early detection gateways for congestion avoidance,” *IEEE/ACM Trans. Networking*, vol. 1, pp. 397–413, Aug. 1993.
- [14] S. Floyd, “TCP and explicit congestion notification,” *ACM Comput. Commun. Rev.*, vol. 24, no. 5, pp. 10–23, Oct. 1994.
- [15] [Online]. Available: <http://www.aciri.org/floyd/REDparameters.txt>
- [16] C. Kenjiro, “A framework for alternate queueing: Toward traffic management by PC-UNIX based routers,” in *Proc. USENIX Annu. Tech. Conf.*, New Orleans, LA, June 1998, pp. 247–258.
- [17] D. Lin and R. Morris, “Dynamics of random early detection,” in *Proc. ACM SIGCOMM*, Sept. 1997, pp. 127–138.
- [18] B. Mah, “An empirical model of HTTP network traffic,” in *Proc. IEEE INFOCOM*, vol. 2, Apr. 1997, pp. 592–600.
- [19] M. May, J. Bolot, C. Diot, and B. Lyles, “Reasons not to deploy RED,” in *Proc. IWQoS*, London, U.K., Mar. 1999, p. 260262.
- [20] M. May, T. Bonald, and J. Bolot, “Analytic evaluation of RED performance,” in *Proc. IEEE INFOCOM*, vol. 3, Mar. 2000, pp. 1415–1424.
- [21] [Online]. Available: <http://www.netstat.net/>
- [22] H. Nielsen, J. Gettys, A. Baird-Smith, E. Prud’hommeaux, H. Lie, and C. Lilley, “Network performance effects of HTTP/1.1, CSS1, and PNG,” in *Proc. ACM SIGCOMM*, Sept. 1997, pp. 155–166.
- [23] T. Ott, T. Lakshman, and L. Wong, “SRED: Stabilized RED,” in *Proc. IEEE INFOCOM*, Mar. 1999, pp. 1346–1355.
- [24] (1998). [Online]. Available: <http://null0.qual.net/brad/papers/reddraft.htm>
- [25] K. Thompson, G. Miller, and R. Wilder, “Wide-area Internet traffic patterns and characteristics,” *IEEE Network*, vol. 11, pp. 10–23, Nov./Dec. 1997.
- [26] C. Villamizar and C. Song, “High-performance TCP in ANSNET,” *ACM Comput. Commun. Rev.*, vol. 24, no. 5, Oct. 1994.

**Mikkel Christiansen**, photograph and biography not available at the time of publication.

**Kevin Jeffay**, photograph and biography not available at the time of publication.

**David Ott**, photograph and biography not available at the time of publication.

**F. Donelson Smith**, photograph and biography not available at the time of publication.