

Internet Engineering Task Force  
INTERNET-DRAFT  
Intended status: Informational  
Expires: August 2008

S. Floyd  
E. Kohler  
Editors  
23 February 2008

**Tools for the Evaluation of Simulation and Testbed Scenarios**  
**draft-irtf-tmrg-tools-05.txt**

Status of this Memo

By submitting this Internet-Draft, each author represents that any applicable patent or other IPR claims of which he or she is aware have been or will be disclosed, and any of which he or she becomes aware will be disclosed, in accordance with [Section 6 of BCP 79](#).

Internet-Drafts are working documents of the Internet Engineering Task Force (IETF), its areas, and its working groups. Note that other groups may also distribute working documents as Internet-Drafts.

Internet-Drafts are draft documents valid for a maximum of six months and may be updated, replaced, or obsoleted by other documents at any time. It is inappropriate to use Internet-Drafts as reference material or to cite them other than as "work in progress."

The list of current Internet-Drafts can be accessed at <http://www.ietf.org/ietf/1id-abstracts.txt>.

The list of Internet-Draft Shadow Directories can be accessed at <http://www.ietf.org/shadow.html>.

This Internet-Draft will expire on August 2008.

Abstract

This document describes tools for the evaluation of simulation and testbed scenarios used in research on Internet congestion control mechanisms. We believe that research in congestion control mechanisms has been seriously hampered by the lack of good models underpinning analysis, simulation, and testbed experiments, and that tools for the evaluation of simulation and testbed scenarios can help in the construction of better scenarios, based on better underlying models. One use of the tools described in this document

is in comparing key characteristics of test scenarios with known characteristics from the diverse and ever-changing real world. Tools characterizing the aggregate traffic on a link include the distribution of per-packet round-trip times, the distribution of connection sizes, and the like. Tools characterizing end-to-end paths include drop rates as a function of packet size and of burst size, the synchronization ratio between two end-to-end TCP flows, and the like. For each characteristic, we describe what aspects of the scenario determine this characteristic, how the characteristic can affect the results of simulations and experiments for the evaluation of congestion control mechanisms, and what is known about this characteristic in the real world. We also explain why the use of such tools can add considerable power to our understanding and evaluation of simulation and testbed scenarios.



Table of Contents

- [1. Introduction. . . . .](#) [5](#)
- [2. Tools . . . . .](#) [6](#)
  - [2.1. Characterizing Aggregate Traffic on a Link . . . . .](#) [6](#)
  - [2.2. Characterizing an End-to-End Path. . . . .](#) [6](#)
  - [2.3. Other Characteristics. . . . .](#) [7](#)
- [3. The Distribution of Per-packet Round-trip Times . . . . .](#) [7](#)
- [4. The Distribution of Connection Sizes. . . . .](#) [8](#)
- [5. The Distribution of Packet Sizes. . . . .](#) [10](#)
- [6. The Ratio Between Forward-path and Reverse-path Traffic. . . . .](#) [10](#)
- [7. The Distribution of Per-Packet Peak Flow Rates. . . . .](#) [11](#)
- [8. The Distribution of Transport Protocols.. . . . .](#) [12](#)
- [9. The Synchronization Ratio . . . . .](#) [12](#)
- [10. Drop or Mark Rates as a Function of Packet Size. . . . .](#) [14](#)
- [11. Drop Rates as a Function of Burst Size.. . . . .](#) [16](#)
- [12. Drop Rates as a Function of Sending Rate.. . . . .](#) [18](#)
- [13. Congestion Control Mechanisms for Traffic, along with Sender and Receiver Buffer Sizes. . . . .](#) [19](#)
- [14. Characterization of Congested Links in Terms of Bandwidth and Typical Levels of Congestion . . . . .](#) [19](#)
  - [14.1. Bandwidth . . . . .](#) [19](#)
  - [14.2. Queue Management Mechanisms . . . . .](#) [19](#)
  - [14.3. Typical Levels of Congestion. . . . .](#) [19](#)
- [15. Characterization of Challenging Lower Layers.. . . . .](#) [19](#)
  - [15.1. Error Losses. . . . .](#) [20](#)
  - [15.2. Packet Reordering . . . . .](#) [20](#)
  - [15.3. Delay Variation . . . . .](#) [21](#)
  - [15.4. Bandwidth Variation . . . . .](#) [22](#)
  - [15.5. Bandwidth and Latency Asymmetry . . . . .](#) [23](#)
  - [15.6. Queue Management Mechanisms . . . . .](#) [24](#)
- [16. Network Changes Affecting Congestion . . . . .](#) [24](#)
  - [16.1. Routing Changes: Routing Loops . . . . .](#) [24](#)
  - [16.2. Routing Changes: Fluttering. . . . .](#) [25](#)
  - [16.3. Routing Changes: Routing Asymmetry . . . . .](#) [26](#)
  - [16.4. Link Disconnections and Intermittent Link Connectivity . . . . .](#) [26](#)
  - [16.5. Changes in Wireless Links: Mobility . . . . .](#) [27](#)
- [17. Using the Tools Presented in this Document . . . . .](#) [27](#)
- [18. Related Work . . . . .](#) [27](#)
- [19. Conclusions. . . . .](#) [27](#)
- [20. Security Considerations. . . . .](#) [27](#)
- [21. IANA Considerations. . . . .](#) [28](#)
- [22. Acknowledgements . . . . .](#) [28](#)
- [Informative References . . . . .](#) [28](#)
- [Authors' Addresses . . . . .](#) [32](#)
- [Full Copyright Statement . . . . .](#) [32](#)

Floyd, Kohler

Expires: August 2008

[Page 3]

Intellectual Property. . . . . [32](#)

TO BE DELETED BY THE RFC EDITOR UPON PUBLICATION:

Changes from [draft-irtf-tmrg-tools-04.txt](#):

- \* Added to the section on "Congestion Control Mechanisms for Traffic". From a contribution from Sara Landstrom.

Changes from [draft-irtf-tmrg-tools-03.txt](#):

- \* No changes.

Changes from [draft-irtf-tmrg-tools-02.txt](#):

- \* Added sections on Challenging Lower Layers and Network Changes affecting Congestion. Contributed by Jasani Rohan, with Julie Tarr, Tony Desimone, Christou Christos, and Vemulapalli Archana.
- \* Minor editing.

Changes from [draft-irtf-tmrg-tools-01.txt](#):

- \* Added section on "Drop Rates as a Function of Sending Rate."
- \* Added a number of new references.

END OF SECTION TO BE DELETED.

## **1. Introduction**

This document discusses tools for the evaluation of simulation and testbed scenarios used in research on Internet congestion control mechanisms. These tools include but are not limited to measurement tools; the tools discussed in this document are largely ways of characterizing aggregate traffic on a link, or characterizing the end-to-end path. One use of these tools is for understanding key characteristics of test scenarios; many characteristics, such as the distribution of per-packet round-trip times on the link, don't come from a single input parameter but are determined by a range of inputs. A second use of the tools is to compare key characteristics of test scenarios with what is known of the same characteristics of the past and current Internet, and with what can be conjectured about these characteristics of future networks. This paper follows the general approach from "Internet Research Needs Better Models" [[FK02](#)].





As an example of the power of tools for characterizing scenarios, a great deal is known about the distribution of connection sizes on a link, or equivalently, the distribution of per-packet sequence numbers. It has been conjectured that a heavy-tailed distribution of connection sizes is an invariant feature of Internet traffic. A test scenario with mostly long-lived traffic, or with a mix with only long-lived and very short flows, does not have a realistic distribution of connection sizes, and can give unrealistic results in simulations or experiments evaluating congestion control mechanisms. For instance, the distribution of connection sizes makes clear the fraction of traffic on a link from medium-sized connections, e.g., with packet sequence numbers from 100 to 1000. These medium-sized connections can slow-start up to a large congestion window, possibly coming to an abrupt stop soon afterwards, contributing significantly to the burstiness of the aggregate traffic, and to the problems facing congestion control.

In the sections below we will discuss a number of tools for describing and evaluating scenarios, show how these characteristics can affect the results of research on congestion control mechanisms, and summarize what is known about these characteristics in real-world networks.

## **2. Tools**

The tools or characteristics that we discuss are the following.

### **2.1. Characterizing Aggregate Traffic on a Link**

- o Distribution of per-packet round-trip times.
- o Distribution of connection sizes.
- o Distribution of packet sizes.
- o Ratio between forward-path and reverse-path traffic.
- o Distribution of peak flow rates.
- o Distribution of transport protocols.

### **2.2. Characterizing an End-to-End Path**

- o Synchronization ratio.
- o Drop rates as a function of packet size.



- o Drop rates as a function of burst size.
- o Drop rates as a function of sending rate.
- o Degree of packet drops.
- o Range of queueing delay.

### **2.3. Other Characteristics**

- o Congestion control mechanisms for traffic, along with sender and receiver buffer sizes.
- o Characterization of congested links in terms of bandwidth and typical levels of congestion (in terms of packet drop rates).
- o Characterization of congested links in terms of buffer size.
- o Characterization of challenging lower layers in terms of reordering, delay variation, packet corruption, and the like.
- o Characterization of network changes affecting congestion, such as routing changes or link outages.

Below we will discuss each characteristic in turn, giving the definition, the factors determining that characteristic, the effect on congestion control metrics, and what is known so far from measurement studies in the Internet.

### **3. The Distribution of Per-packet Round-trip Times**

Definition: The distribution of per-packet round-trip times on a link is defined formally by assigning to each packet the most recent round trip time measured for that end-to-end connection. In practice, coarse-grained information is generally sufficient, even though it has been shown that there is significant variability in round-trip times within a TCP connection [[AKSJ03](#)], and it is sufficient to assign to each packet the first round-trip time measurement for that connection, or to assign the current round-trip time estimate maintained by the TCP connection.

Determining factors: The distribution of per-packet round-trip times on a link is determined by end-to-end propagation delays, by queueing delays along end-to-end paths, and by the congestion control mechanisms used by the traffic. For example, for a scenario using TCP, TCP connections with smaller round-trip times will receive a proportionally larger fraction of traffic than competing



TCP connections with larger round-trip times, all else being equal, due to the dynamics of TCP favoring flows with smaller round-trip times. This will generally shift the distribution of per-packet RTTs lower relative to the distribution of per-connection RTTs, since short-RTT connections will have more packets.

Effect on congestion control metrics: The distribution of per-packet round-trip times on a link affects the burstiness of the aggregate traffic, and therefore can affect congestion control performance in a range of areas such as delay/throughput tradeoffs. The distribution of per-packet round-trip times can also affect metrics of fairness, degree of oscillations, and the like. For example, long-term oscillations of queueing delay are more likely to occur in scenarios with a narrow range of round-trip times [[FK02](#)].

Measurements: The distribution of per-packet round-trip times for TCP traffic on a link can be measured from a packet trace with the passive TCP round-trip time estimator from Jiang and Dovrolis [[JD02](#)]. [Add pointers to other estimators, such as ones mentioned in JD02. Add a pointer to Mark Allman's loss detection tool.] Their paper shows the distribution of per-packet round-trip times for TCP packets for a number of different links. For the links measured, the percent of packets with round-trip times at most 100 ms ranged from 30% to 80%, and the percent of packets with round-trip times at most 200 ms ranged from 55% to 90%, depending on the link.

In the NS simulator, the distribution of per-packet round-trip times for TCP packets on a link can be reported by the queue monitor, using TCP's estimated round-trip time added to packet headers. This is illustrated in the validation test `./test-all-simple stats3` in the directory `tcl/test`.

Scenarios: [[FK02](#)] shows a relatively simple scenario, with a dumbbell topology with four access links on each end, that gives a fairly realistic range of round-trip times. [Look for the other citations to add.]

#### **4. The Distribution of Connection Sizes**

Definition: Instead of the connection-based measurement of the distribution of connection sizes (the total number of bytes or of data packets in a connection), we consider the packet-based measurement of the distribution of packet sequence numbers. The distribution of packet sequence numbers on a link is defined by giving each packet a sequence number, where the first packet in a connection has sequence number 1, the second packet has sequence number 2, and so on. The distribution of packet sequence numbers



can be derived in a straightforward manner from the distribution of connection sizes, and vice versa; however, the distribution of connection sizes is more suited for traffic generators, and the distribution of packet sequence numbers is more suited for measuring and illustrating the packets actually seen on a link over a fixed interval of time. There has been a considerably body of research over the last ten years on the heavy-tailed distribution of connection sizes for traffic on the Internet. [[CBC95](#)] [Add citations.]

Determining factors: The distribution of connection sizes is largely determined by the traffic generators used in a scenario. For example, is there a single traffic generator characterized by a distribution of connection sizes? A mix of long-lived and web traffic, with the web traffic characterized by a distribution of connection sizes? Or something else?

Effect on congestion control metrics: The distribution of packet sequence numbers affects the burstiness of aggregate traffic on a link, thereby affecting all congestion control metrics for which this is a factor. As an example, [[FK02](#)] illustrates that the traffic mix can affect the queue dynamics on a congested link. [Find more to cite, about the effect of the distribution of packet sequence numbers on congestion control metrics.]

[Add a paragraph about the impact of medium-size flows.]

[Add a paragraph about the impact of flows starting and stopping.]

[Add a warning about scenarios that use only long-lived flows, or a mix of long-lived and very short flows.]

Measurements: [Cite some of the literature.]

Traffic generators: Some of the available traffic generators are listed on the web site for "Traffic Generators for Internet Traffic" [[TG](#)]. This includes pointers to traffic generators for peer-to-peer traffic, traffic from online games, and traffic from Distributed Denial of Service (DDoS) attacks.

In the NS simulator, the distribution of packet sequence numbers for TCP packets on a link can be reported by the queue monitor at a router. This is illustrated in the validation test `./test-all-simple stats3` in the directory `tcl/test`.





## **5. The Distribution of Packet Sizes**

Definition: The distribution of packet sizes is defined in a straightforward way, using packet sizes in bytes.

Determining factors: The distribution of packet sizes is determined by the traffic mix, the path MTUs, and by the packet sizes used by the transport-level senders.

The distribution of packet sizes on a link is also determined by the mix of forward-path TCP traffic and reverse-path TCP traffic in that scenario, for a scenario characterized by a `forward path' (e.g., left to right on a particular link) and a `reverse path' (e.g., right to left on the same link). For such a scenario, the forward-path TCP traffic contributes data packets to the forward link and acknowledgment packets to the reverse link, while the reverse-path TCP traffic contributes small acknowledgment packets to the forward link. The ratio between TCP data and TCP ACK packets on a link can be used as some indication of the ratio between forward-path and reverse-path TCP traffic.

Effect on congestion control metrics: The distribution of packet sizes on a link is an indicator of the ratio of forward-path and reverse-path TCP traffic in that network. The amount of reverse-path traffic determines the loss and queueing delay experienced by acknowledgement packets on the reverse path, significantly affecting the burstiness of the aggregate traffic on the forward path. [In what other ways does the distribution of packet sizes affect congestion control metrics?]

Measurements: There has been a wealth of measurements over time on the packet size distribution of traffic [[A00](#)], [[HMTG01](#)]. These measurements are generally consistent with a model of roughly 10% of the TCP connections using an MSS of roughly 500 bytes, and with the other 90% of TCP connections using an MSS of 1460 bytes.

## **6. The Ratio Between Forward-path and Reverse-path Traffic**

Definition: For a scenario characterized by a `forward path' (e.g., left to right on a particular link) and a `reverse path' (e.g., right to left on the same link), the ratio between forward-path and reverse-path traffic can be defined as the ratio between the forward-path traffic in bps, and the reverse-path traffic in bps.

Determining factors: The ratio between forward-path and reverse-path traffic is defined largely by the traffic mix.



Effect on congestion control metrics: Zhang, Shenker and Clark have shown in 1991 that for TCP, the amount of reverse-path traffic affects the ACK compression and packet drop rate for TCP acknowledgement packets, significantly affecting the burstiness of TCP traffic on the forward path [[ZSC91](#)]. The queueing delay on the reverse path also affects the performance of delay-based congestion control mechanisms, if the delay is computed based on round-trip times. This has been shown by Grieco and Mascolo in [[GM04](#)] and by Prasad, Jain, and Dovrolis in [[PJD04](#)].

Measurements: There is a need for measurements on the range of ratios between forward-path and reverse-path traffic for congested links. In particular, for TCP traffic traversing congested link X, what is the likelihood that the acknowledgement traffic will encounter congestion (i.e., queueing delay, packet drops) somewhere on the reverse path as well?

As discussed in [Section 5](#), the distribution of packet sizes on a link can be used as an indicator of the ratio of forward-path and reverse-path TCP traffic in that network.

## **7. The Distribution of Per-Packet Peak Flow Rates**

Definition: The distribution of peak flow rates is defined by assigning to each packet the peak sending rate in bytes per second of that connection, where the peak sending rate is defined over 0.1-second intervals. The distribution of peak flow rates gives some indication of the ratio of "alpha" and "beta" traffic on a link, where alpha traffic on a congested link is defined as traffic with that link at the main bottleneck, while the beta traffic on the link has a primary bottleneck elsewhere along its path [[RSB01](#)].

Determining factors: The distribution of peak flow rates is determined by flows with bottlenecks elsewhere along their end-to-end path, e.g., flows with low-bandwidth access links. The distribution of peak flow rates is also affected by applications with limited sending rates.

Effect on congestion control metrics: The distribution of peak flow rates affects the burstiness of aggregate traffic, with low-peak-rate traffic decreasing the aggregate burstiness, and adding to the traffic's tractability.

Measurements: [[RSB01](#)]. The distribution of peak rates can be expected to change over time, as there is an increasing number of high-bandwidth access links to the home, and of high-bandwidth Ethernet links at work and at other institutions.



Simulators: [For NS, add a pointer to the DelayBox, "http://dirt.cs.unc.edu/delaybox/", for more easily simulating low-bandwidth access links for flows.]

Testbeds: In testbeds, Dummynet [[Dummynet](#)] and NISTNet [[NISTNet](#)] provide convenient ways to emulate paths with different limited peak rates.

## 8. The Distribution of Transport Protocols.

Definition: The distribution of transport protocols on a congested link is straightforward, with each packet given its associated transport protocol (e.g., TCP, UDP). The distribution is often given both in terms of packets and in terms of bytes.

For UDP packets, it might be more helpful to classify them in terms of the port number, or the assumed application (e.g., DNS, RIP, games, Windows Media, RealAudio, RealVideo, etc.) [[MAWI](#)]. Other traffic includes ICMP, IPSEC, and the like. In the future there could be traffic from SCTP, DCCP, or from other transport protocols.

Effect on congestion control metrics: The distribution of transport protocols affects metrics relating to the effectiveness of AQM mechanisms on a link.

Measurements: In the past, TCP traffic has typically consisted of 90% to 95% of the bytes on a link [[UW02](#)], [[UA01](#)]. [Get updated citations for this.] Measurement studies show that TCP traffic from web servers almost always uses conformant TCP congestion control procedures [[MAF05](#)].

## 9. The Synchronization Ratio

Definition: The synchronization ratio is defined as the degree of synchronization of loss events between two TCP flows on the same path. Thus, the synchronization ratio is defined as a characteristic of an end-to-end path. When one TCP flow of a pair has a loss event, the synchronization ratio is given by the fraction of those loss events for which the second flow has a loss event within one round-trip time. Each connection in a flow pair has a separate synchronization ratio, and the overall synchronization ratio of the pair of flows is the higher of the two ratios. When measuring the synchronization ratio, it is preferable to start the two TCP flows at slightly different times, with large receive windows.

Determining factors: The synchronization ratio is determined largely by the traffic mix on the congested link, and by the AQM mechanism



(or lack of AQM mechanism).

Different types of TCP flows are also likely to have different synchronization measures. E.g., Two HighSpeed TCP flows might have higher synchronization measures than two Standard TCP flows on the same path, because of their more aggressive window increase rates. Raina, Towsley, and Wischik [[RTW05](#)] have discussed the relationships between synchronization and TCP's increase and decrease parameters.

Effect on congestion control metrics: The synchronization ratio affects convergence times for high-bandwidth TCPs. Convergence times are known to be poor for some high-bandwidth protocols in environments with high levels of synchronization [[LS06](#)]. However, the scenarios in [[LS06](#)] are of a congested link with one-way traffic, long-lived flows all with the same round-trip time, and Drop-Tail queue management at routers. These are not realistic scenarios; instead, these are the scenarios that I assume would maximize the degree of synchronization between flows.

Wischik and McKeown [[WM05](#)] have shown that the level of synchronization affects the buffer requirements at congested routers. Baccelli and Hong [[BH02](#)] have a model showing the effect of the synchronization ratio on aggregate throughput.

Measurements: Grenville Armitage and Qiang Fu have performed initial experiments of synchronization in the Internet, using Standard TCP flows, and have found very low levels of synchronization.

In a discussion of the relationship between stability and desynchronization, Raina, Towsley, and Wischik [[RTW05](#)] report that "synchronization has been reported again and again in simulations". In contrast, synchronization has not been reported again and again in the real-world Internet.

Appenzeller, Keslassy, and McKeown in [[AKM04](#)] report the following: "Flows are not synchronized in a backbone router carrying thousands of flows with varying RTTs. Small variations in RTT or processing time are sufficient to prevent synchronization [[QZK01](#)]; and the absence of synchronization has been demonstrated in real networks [[F02](#), [IMD01](#)]."

[Appenzeller et al, Sizing Router Buffers, reports that synchronization is rare as the number of competing flows increases. Kevin Jeffay has some results on synchronization also.]

Needed: We need measurements of the synchronization ratio for flows that use high-bandwidth protocols over high-bandwidth paths, given typical levels of competing traffic and with typical queueing





mechanisms at routers (whatever these are), to see if there are higher levels of synchronization with high-bandwidth protocols such as HighSpeed TCP, Fast TCP, and the like, which are more aggressive than Standard TCP. The assumption would be that in many environments, high-bandwidth protocols have higher levels of synchronization than flows using Standard TCP.

#### **10. Drop or Mark Rates as a Function of Packet Size**

Definition: Drop rates as a function of packet size are defined by the actual drop rates for different packets on an end-to-end path or on a congested link over a particular time interval. In some cases, e.g., Drop-Tail queues in units of packets, general statements can be made; e.g., that large and small packets will experience the same packet drop rates. However, in other cases, e.g., Drop-Tail queues in units of bytes, no such general statement can be made, and the drop rate as a function of packet size will be determined in part by the traffic mix at the congested link at that point of time.

Determining factors: The drop rate as a function of packet size is determined in part by the queue architecture. E.g., is the Drop-Tail queue in units of packets, of bytes, of 60-byte buffers, or of a mix of buffer sizes? Is the AQM mechanism in packet mode, dropping each packet with the same probability, or in byte mode, with the probability of dropping or marking a packet being proportional to the packet size in bytes.

The effect of packet size on drop rate would also be affected by the presence of preferential scheduling for small packets, or by differential scheduling for packets from different flows (e.g., per-flow scheduling, or differential scheduling for UDP and TCP traffic).

In many environments, the drop rate as a function of packet size will be heavily affected by the traffic mix at a particular time. For example, is the traffic mix dominated by large packets, or by smaller ones? In some cases, the overall packet drop rate could also affect the relative drop rates for different packet sizes.

In wireless networks, the drop rate as a function of packet size is also determined by the packet corruption rate as a function of packet size. [Cite Deborah Pinck's papers on Satellite-Enhanced Personal Communications Experiments and on Experimental Results from Internetworking Data Applications Over Various Wireless Networks Using a Single Flexible Error Control Protocol.] [Cite the general literature.]



Effect on congestion control metrics: The drop rate as a function of packet size has a significant effect on the performance of congestion control for VoIP and other small-packet flows.

[Citation: "TFRC for Voice: the VoIP Variant", [draft-ietf-dccp-tfrc-voip-02.txt](#), and earlier papers.] The drop rate as a function of packet size also has an effect on TCP performance, as it affects the drop rates for TCP's SYN and ACK packets. [Citation: Jeffay and others.]

Measurements: We need measurements of the drop rate as a function of packet size over a wide range of paths, or for a wide range of congested links. For tests of relative drop rates on end-to-end packets, one possibility would be to run successive TCP connections with 200-byte, 512-byte, and 1460-byte packets, and to compare the packet drop rates. The ideal test would include running TCP connections on the reverse path, to measure the drop rates for the small ACK packets on the forward path. It would also be useful to characterize the difference in drop rates for 200-byte TCP packets and 200-byte UDP packets, even though some of this difference could be due to the relative burstiness of the different connections.

Ping experiments could also be used to get measurements of drop rates as a function size, but it would be necessary to make sure that the ping sending rates were adjusted to be TCP-friendly.

[Cite the known literature on drop rates as a function of packet size.]

Our conjecture is that there is a wide range of behaviors for this characteristic in the real world. Routers include Drop-Tail queues in packets, bytes, and buffer sizes in between; these will have quite different drop rates as a function of packet size. Some routers include RED in byte mode (the default for RED in Linux) and some have RED in packet mode (Cisco, I believe). This also affects drop rates as a function of packet size.

Some routers on congested access links use per-flow scheduling. In this case, does the per-flow scheduling have the goal of fairness in \*bytes\* per second or in \*packets\* per second? What effect does the per-flow scheduling have on the drop rate as a function of packet size, for packets in different flows (e.g., a small-packet VoIP flow competing against a large-packet TCP flow) or for packets within the same flow (small ACK packets and large data packets on a two-way TCP connection).



## **11. Drop Rates as a Function of Burst Size.**

Definition: Burst-tolerance, or drop rates as a function of burst size, can be defined in terms of an end-to-end path, or in terms of aggregate traffic on a congested link.

The burst-tolerance of an end-to-end path is defined in terms of connections with different degrees of burstiness within a round-trip time. When packets are sent in bursts of  $N$  packets, does the drop rate vary as a function of  $N$ ? For example, if the TCP sender sends small bursts of  $K$  packets, for  $K$  less than the congestion window, how does the size of  $K$  affect the loss rate? Similarly, for a ping tool sending pings at a certain rate in packets per second, one could see how the clustering of the ping packets in clusters of size  $K$  affects the packet drop rate. As always with such ping experiments, it would be important to adjust the sending rate to maintain a longer-term sending rate that was TCP-friendly.

Determining factors: The burst-tolerance is determined largely by the AQM mechanisms for the congested routers on a path, and by the traffic mix. For a Drop-Tail queue with only a small number of competing flows, the burst-tolerance is likely to be low, and for AQM mechanisms where the packet drop rate is a function of the average queue size rather than the instantaneous queue size, the burst tolerance should be quite high.

Effect on congestion control metrics: The burst-tolerance of the path or congested link can affect fairness between competing flows with different round-trip times; for example, Standard TCP flows with longer round-trip times are likely to have a more bursty arrival pattern at the congested link than that of Standard TCP flows with shorter round-trip times. As a result, in environment with low burst tolerance (e.g., scenarios with Drop-Tail queues), longer-round-trip-time TCP connections can see higher packet drop rates than other TCP connections, and receive an even smaller fraction of the link bandwidth than they would otherwise. [FJ92] (Section 3.2). We note that some TCP traffic is inherently bursty, e.g., Standard TCP without rate-based pacing, particularly in the presence of dropped ACK packets or of ACK compression. The burst-tolerance of a router can also affect the delay-throughput tradeoffs and packet drop rates of the path or of the congested link.

Measurements: One could measure the burst-tolerance of an end-to-end path by running successive TCP connections, forcing bursts of size at least  $K$  by dropping an appropriate fraction of the ACK packets to the TCP receiver. Alternately, if one had control of the TCP sender, one could modify the TCP sender to send bursts of  $K$  packets when the congestion window was  $K$  or more segments.

Floyd, Kohler

Expires: August 2008

[Section 11](#). [Page 16]

Blanton and Allman in [[BA05](#)] consider the TCP micro-bursts that result from the receipt of a single acknowledgement packet or from application-layer dynamics, and consider bursts of four or more packets. They consider four traces, and plot the probability of at least one packet from a burst being lost, as a function of burst size. Considering only connections with both bursts and packet losses, the probability of packet loss when the TCP connection was bursting was somewhat higher than the probability of packet loss when the TCP connection was not bursting in three of the four traces. For each trace, the paper shows the aggregate probability of loss as a function of the burst size in packets. Because these are aggregate statistics, it cannot be determined if there is a correlation between the burst size and the TCP connection's sending rate.

[Look at: M. Allman and E. Blanton, "Notes on Burst Mitigation for Transport Protocols", ACM Computer Communication Review, vol. 35(2), (2005).]

Making inferences about the AQM mechanism for the congested router on an end-to-end path: One potential use of measurement tools for determining the burst-tolerance of an end-to-end path would be to make inferences about the presence or absence of an AQM mechanism at the congested link or links. As a simple test, one could run a TCP connection until the connection comes out of slow-start. If the receive window of the TCP connection was sufficiently high that the connection exited slow-start with packet drops or marks instead of because of the limitation of the receive window, one could record the congestion window at the end of slow-start, and the number of packets dropped from this window. A high packet drop rate might be more typical of a Drop-Tail queue with small-scale statistical multiplexing on the congested link, and a single packet drop coming out of slow-start would suggest an AQM mechanism at the congested link.

The synchronization measure could also add information about the likely presence or absence of AQM on the congested link(s) of an end-to-end path, with paths with higher levels of synchronization being more likely to have Drop-Tail queues with small-scale statistical multiplexing on the congested link(s).

Lui and Crovella in [[LC01](#)] use loss pairs to infer the queue size when packets are dropped. A loss pair consists of two packets sent back-to-back, where one of the two packets is dropped in the network. The round-trip time of the surviving packet is used to estimate the round-trip time when the companion packet was dropped in the network. For a path with Drop-Tail queueing at the congested link, this round-trip time can be used to estimate the queue size,

Floyd, Kohler

Expires: August 2008

[Section 11](#). [Page 17]



given estimates of the link bandwidth and minimum round-trip time. For a path with AQM at the congested link, trial pairs are also considered, where a trial pair is any pair of packets sent back-to-back. [LC01] uses the ratio between the number of loss pairs and the number of trial pairs for each round-trip range to estimate the drop probability of the AQM mechanism at the congested link as a function of queue size. [LC01] uses loss pairs in simulation settings with a minimum of noise in terms of queueing delays elsewhere on the forward or reverse path.

[Cite the relevant literature about tools for determining the AQM mechanism on an end-to-end path.]

## **12. Drop Rates as a Function of Sending Rate.**

Definition: Drop rates as a function of sending rate is defined in terms of the drop behavior of a flow in the end-to-end path. That is, does the sending rate of an individual flow affect its own packet drop rate, or its packet drop rate largely independent of the sending rate of the flow?

Determining factors: The sending rate of the flow affects its own packet drop rate in an environment with small-scale statistical multiplexing on the congested link. The packet drop rate is largely independent of the sending rate in an environment with large-scale statistical multiplexing, with many competing small flows at the congested link. Thus, the behavior of drop rates as a function of sending rate is a rough measure of the level of statistical multiplexing on the congested links of an end-to-end path.

Effect on congestion control metrics: The level of statistical multiplexing at the congested link can affect the performance of congestion control mechanisms in transport protocols. For example, delay-based congestion control is often better suited to environments small-scale statistical multiplexing at the congested link, where the transport protocol responds to the delay caused by its own sending rate.

Measurements: In a simulation or testbed, the level of statistical multiplexing on the congested link can be observed directly. In the Internet, the level of statistical multiplexing on the congested links of an end-to-end path can be inferred indirectly through per-flow measurements, by observing whether the packet drop rate varies as a function of the sending rate of the flow.



### **13. Congestion Control Mechanisms for Traffic, along with Sender and Receiver Buffer Sizes.**

Effect on congestion control metrics: Please don't evaluate AQM mechanisms by using Reno TCP, or evaluate new transport protocols by comparing them with the performance of Reno TCP. For measurement data, see below. For a more detailed explanation, see [[FK02](#)] ([Section 3.4](#)).

SACK and DSACK: Medina et al. in [[MAF05](#)] tested 84,394 servers for SACK capability. Of these, the majority, 68%, were SACK-Capable. Approximately half of the SACK-Capable web servers supported DSACK.

Allman in [[A00](#)] reports that the percentage of web clients that were SACK-Capable increased from 8% in December 1998 to 40% in March 2000. This trend continued, with 88% of the clients advertising 'SACK permitted' in the 2004 data reported in [[MAF05](#)]. Only 3% of the clients sent DSACKs, but this number does not reveal how many clients would have sent DSACKs upon receiving duplicate data.

Reno and NewReno: When the TBIT client used by Medina et al. in [[MAF05](#)] pretended not to be SACK-Capable, only 33% of the web servers were classified as NewReno, Reno, Tahoe, or Other, but of these, the majority (76%) were classified as NewReno, and 15% were classified as Reno. In [[PF01](#)] NewReno was already observed as the dominating congestion control algorithm in the absence of SACK information. Out of 3,728 web servers, 1,571 performed NewReno congestion control in that investigation.

### **14. Characterization of Congested Links in Terms of Bandwidth and Typical Levels of Congestion**

#### **14.1. Bandwidth**

#### **14.2. Queue Management Mechanisms**

#### **14.3. Typical Levels of Congestion**

[Pointers to the current state of our knowledge.]

### **15. Characterization of Challenging Lower Layers.**

With an increasing number of wireless networks connecting to the wired Internet, more and more end-to-end paths will contain a combination of wired and wireless links. These wireless links



exhibit new characteristics which congestion control mechanisms will need to cope with. The main characteristics, detailed in subsequent sections, include error losses, packet reordering, delay variation, bandwidth variation, and bandwidth and latency asymmetry.

### **15.1. Error Losses**

**Definition:** Packet losses due to corruption rarely occur on wired links, but occur on wireless links due to random/transient errors and/or extended burst errors. If packet errors cannot be detected and discarded within the network through error detection schemes or recovered through error recovery schemes such as Forward Error Correction (FEC) and Automatic Repeat Request (ARQ), the corrupted packet is discarded, resulting in an error loss.

**Determining Factors:** Error losses are primarily caused by the degradation of the quality of a wireless link (multipath, fade, etc.). Link errors can be characterized by the type of errors that occur (e.g., random, burst), the length of time they occur, and the frequency at which they occur. These characteristics are highly dependent on the wireless channel conditions and are influenced by the distance between two nodes on a wireless link, the type and orientation of antennas, encoding algorithms, and other factors [22]. Therefore, error losses are significantly influenced by these link errors.

**Effect on congestion control metrics:** Since error losses can be unrelated to congestion, congestion control mechanisms should recover from these types of losses differently than from congestion losses. If congestion control mechanisms misinterpret error losses as congestion losses, then they respond inappropriately, reducing the sending rate too much [23]. As a result, an unnecessary reduction in the sending rate can occur, when in reality the available bandwidth has not changed. This can result in a reduction in throughput and underutilization of the channel. However, error recovery mechanisms such as FEC or ARQ are heavily used in cellular networks to reduce the impact of error losses [IMLGK03].

**Measurements:** In 3G cellular networks, error recovery mechanisms have reduced the rate of error losses to under 1%, making their impact marginal [CR04].

### **15.2. Packet Reordering**

**Definition:** Due to the connectionless nature of IP, packets can arrive out of order at their destination. Packet reordering events can occur at varying times and to varying degrees. For example, a particular channel may reorder one out of ten packets and the



reordered packet arrives three packets out of order.

Determining Factors: For the most part, packet reordering on wireless links rarely occurs. However, packet re-ordering can occur due to link layer error recovery. Extensive packet reordering has been shown to occur with particular handoff mechanisms, and is definitely detrimental to transport performance [[GF04](#)].

Effects on congestion control metrics: With TCP, packet reordering can cause the receiver to wait for the arrival of packets that are out of order, since the receiver must reassemble the packets in the correct order before passing them up to the application. With TCP and other transport protocols, packet reordering can also result in the sender incorrectly inferring packet loss, triggering packet retransmissions and congestion control responses.

Measurements: Measurements by Zhou and Miegheem show that reordering happens quite often in the Internet, but few streams have more than two reordered packets [[ZM04](#)]. For further measurements, see [[ANP06](#)][[BPS99](#)][[LC05](#)].

### **15.3. Delay Variation**

Definition: Delay Variation occurs when selected packets of a given flow experience a difference in the One-Way-Delay across a network. Delay variation can be caused by a variation in propagation, transmission and queueing delay that can occur across links or network nodes.

Determining Factors: Delay and delay variation is introduced due to various features of wireless links [[IMLGK03](#)]. The delay experience by subsequent packets of a given flow can change due to On-Demand Resource Allocation, which allocates a wireless channel to a user based on current bandwidth availability. In addition, FEC and ARQ, which are commonly used to combat error loss on wireless links, can introduce delay into the channel, depending on the degree of error loss that occurs. These mechanisms either resend packets that have been corrupted or attempt to recover the actual corrupted packet, which both add delay to the channel.

Effect on congestion control metrics: A spike in delay can have a negative impact on transport protocols [[AAR03](#)][[AGR00](#)][[CR04](#)]. Transport protocols use timers for loss recovery and for congestion control, which are set according to the RTT. Delay spikes can trigger spurious timeouts that cause unnecessary retransmissions and incorrect congestion control responses. If these delay spikes continue, they can inflate the retransmission timeout, increasing the wait before a dropped packet is recovered. Delay-based





congestion control mechanisms (e.g. TCP Vegas, TCP Westwood, etc.) use end-to-end delay to control the sending rate of the sender. Delay-based congestion control mechanisms use delay to indicate when there is congestion in the network. When delay variation occurs for reasons other than queueing delay, delay based congestion control mechanisms can reduce the sending rate unnecessarily. Rate-based protocols can perform poorly as they do not adjust the sending rate after a change in the RTT, possibly creating unnecessary congestion [GF04].

Measurements: Cellular links, particularly GPRS and CDMA2000, can have one-way latencies varying from 100 to 500 ms [IMLGK03]. The length of a delay variation event can vary from three to fifteen seconds and the frequency at which delay variation events occur can be anywhere from 40 to 400 seconds. GEO satellite links tend not see much variation in delay, while LEO satellite links can see significant variability in delay due to the constant motion of satellites and multiple hops. The delay variation of LEO satellite links can be from 40 to 400ms [GK98].

#### **15.4. Bandwidth Variation**

Definition: The bandwidth of a wireless channel can vary over time during a single session, as wireless networks can change the available bandwidth allotted to a user. Therefore, a user may have a low-bandwidth channel for part of their session and a high-bandwidth channel the remainder of the session. The bandwidth of the channel can vary abruptly or gradually, at various intervals, and these variations can occur at different times. On-demand Resource Allocation is one of the mechanisms used to dynamically allocate resources to users according to system load and traffic priority. For instance, in GPRS a radio channel is allocated when data arrives toward the user, and released when the queue size falls below a certain threshold [GPAR02][W01].

Determining Factors: The amount of bandwidth of a given channel that is allocated by the wireless network to a user can vary based on a number of factors. Factors such as wireless conditions or the amount of users connected to the base station both affect the available capacity [IMLGK03]. In certain satellite systems, technologies such as On-Demand Resource Allocation constantly adjust the available bandwidth for a given user every second, which can cause significant bandwidth variation during a session. On-Demand Resource Allocation is designed into most 2.5 and 3G wireless networks, but it can be implemented differently from network to network, resulting in different impacts on the link.

Effect on congestion control metrics: In the absence of congestion,

Floyd, Kohler

Expires: August 2008 [Section 15.4](#). [Page 22]

congestion control mechanisms increase the sending rate gradually over multiple round-trip times. If the bandwidth of a wireless channel suddenly increases and this increases the bandwidth available on the end-to-end path, the transport protocol might not be able to increase its sending rate quickly enough to use the newly-available bandwidth [24]. If the bandwidth of a wireless channel suddenly decreases and this decreases the bandwidth available on the end-to-end path, the sender might not decrease its sending rate quickly enough, resulting in transient congestion. Frequent changes of bandwidth on a wireless channel can result in the average transmission rate of the channel being limited by the amount of bandwidth available during times where the channel has the lowest bandwidth. Persistent delay variation can inflate the retransmission timeout, increasing the wait before a dropped packet is recovered, ultimately leading to channel underutilization [GPAR02][W01].

Measurements: Further references on the measurements for the amount of bandwidth variation are needed. On-demand channel allocation can be modeled by introducing an additional delay when a packet arrives to a queue that has been empty longer than the channel hold time (i.e., propagation delay). The delay value represents the channel allocation delay, and the hold time represents the duration of channel holding after transmitting a data packet [GPAR02][W01]. See also [NM01].

### **15.5. Bandwidth and Latency Asymmetry**

Definition: The bandwidth in the forward direction or uplink can be different than the bandwidth in the reverse direction or downlink. Similar to bandwidth asymmetry, latency in the forward direction or uplink can be different than latency in the reverse direction or downlink. For example, bandwidth asymmetry occurs in wireless networks where the channel from the mobile to base station (uplink) has a fraction of the bandwidth of the channel from the base station to the mobile channel (downlink).

Determining Factors: Bandwidth and latency asymmetry can occur for a variety of reasons. Mobile devices that must transmit at lower power levels to conserve power have low bandwidth and high latency transmission, while base stations can transmit at higher power levels, resulting in higher bandwidth and lower latency [IMLGK03]. In addition, because applications such as HTTP require significantly more bandwidth on the downlink as opposed to the uplink, wireless networks have been designed with asymmetry to accommodate these applications. Coupled with these design constraints, the environmental conditions can add increased asymmetry [HK99].

Floyd, Kohler

Expires: August 2008 [Section 15.5](#). [Page 23]

Effect on congestion control metrics: TCP's congestion control algorithms rely on ACK-clocking, with the reception of ACKs controlling sending rates. If ACKs are dropped or delayed in the reverse direction, then the sending rate in the forward direction can be reduced. In addition, excessive delay can result in a retransmit timeout and a corresponding reduction in the sending rate [[HK99](#)][23].

Measurements: For cellular networks the downlink bandwidth typically does not exceed three to six times the uplink bandwidth [[IMLGK03](#)]. However, different cellular networks (e.g. IS-95, CDMA2000, etc.) have different ratios of bandwidth and latency asymmetry.

### **[15.6.](#) Queue Management Mechanisms**

In wireless networks, queueing delay typically occurs at the end points of a wireless connection (i.e. mobile device, base station) [[GF04](#)].

Measurements: For current cellular and WLAN links, the queue can plausibly be modeled with Drop-Tail queueing with a configurable maximum size in packets. The use of RED may be more appropriate for modeling satellite or future cellular and WLAN links [[GF04](#)].

## **[16.](#) Network Changes Affecting Congestion**

Changes in the network can have a significant impact on the performance of congestion control algorithms. These changes can include events such as the unnecessary duplication of packets, topology changes due to node mobility, and temporary link disconnections. These types of network changes can be broadly categorized as routing changes, link disconnections and intermittent link connectivity, and mobility.

### **[16.1.](#) Routing Changes: Routing Loops**

Definition: A routing loop is a network event in which packets continue to be routed in an endless circle until the packets are eventually dropped [[P96](#)].

Determining factors: Routing loops can occur when the network experiences a change in connectivity which is not immediately propagated to all of the routers [[H95](#)]. Network and node mobility are examples of network events that can cause a change in connectivity.

Effect on Congestion Control: Loops can rapidly lead to congestion,



as a router will route packets towards a destination, but the forwarded packets end up being routed back to the router. Furthermore, network congestion due to a routing loop with multicast packets will be more severe than with unicast packets because each router replicates multicast packets thereby causing congestion more rapidly [P96].

Measurements: Routing dynamics that lead to temporary route loss or forwarding loops are also called routing failures. ICMP response messages, measured by traceroutes and pings, can be used to identify routing failures [WMWGB06].

## **16.2. Routing Changes: Fluttering**

Definition: The term fluttering is used to describe rapidly-oscillating routing. Fluttering occurs when a router alternates between multiple next-hop routers in order to split the load among the links to those routers. While fluttering can provide benefits as a way to balance load in a network, it also creates problems for TCP [P96][AP99].

Determining factors: Multi-path routing in the Internet can cause route fluttering. Route fluttering can result in significant out-of-order packet delivery and/or frequent abrupt end-to-end RTT variation [P97].

Effect on Congestion Control Metrics: When two routes have different propagation delays, packets will often arrive at the destination out-of-order, depending on whether they arrived via the shorter route or the longer route. Whenever a TCP endpoint receives an out-of-order packet, this triggers the transmission of a duplicate acknowledgement to inform the sender that the receiver has a hole in its sequence space. If three out-of-order packets arrive in a row, then the receiver will generate three duplicate acknowledgements for the segment that was not received. These duplicate acknowledgements will trigger fast retransmit by the sender, leading it to reduce the sending rate and needlessly retransmit data. Thus, out-of-order delivery can result in unnecessary reductions in the sending rate and also in redundant network traffic, due to extra acknowledgements and possibly unnecessary data retransmissions [P96][AP99].

Measurements: Two metrics [WMWGB06] can be used to measure the degree of out-of-order delivery: the number of reordering packets and the reordering offset. The number of reordering packets is simply the number of packets that are considered out of order. The reordering offset for an out-of-order packet is the difference between the actual arrival order and the expected arrival order. See also [LMJ96].

Floyd, Kohler

Expires: August 2008 [Section 16.2](#). [Page 25]



### **16.3. Routing Changes: Routing Asymmetry**

Definition: Routing asymmetry occurs when packets traveling between two end-points follow different routes in the forward and reverse directions. The two routes could have different characteristics in terms of bandwidth, delay, levels of congestion, etc. [P96].

Determining factors: Some of the main causes for asymmetry are policy routing, traffic engineering, and the absence of a unique shortest path between a pair of hosts. While the lack of a unique shortest path is one potential contributor to asymmetric routing within domains, the principal source of asymmetries in backbone routers is policy routing. Another cause of routing asymmetry is adaptive routing, in which a router shifts traffic from a highly loaded link to a less loaded one, or load balances across multiple paths [P96].

Effect on Congestion Control: When delay-based congestion control is used, asymmetry can introduce problems in estimating the one-way latency between hosts.

Measurements: Further references are needed.

### **16.4. Link Disconnections and Intermittent Link Connectivity**

Definition: A link disconnection is a period when the link loses all frames, until the link is restored. Intermittent Link Connectivity occurs when the link is disconnected regularly and for short periods of time. This is a common characteristic of wireless links, particularly those with highly mobile nodes [AP99].

Determining factors: In a wireless environment, link disconnections and intermittent link connectivity could occur when a mobile device leaves the range of a base station, which can lead to signal degradation or failure in a handoff [AP99].

Effect on Congestion Control Metrics: If a link disconnection lasts longer than the TCP RT0, and results in a path disconnection for that period of time, the TCP sender will perform a retransmit timeout, resending a packet and reducing the sending rate. TCP will continue this pattern, with longer and longer retransmit timeouts, up to a retry limit, until an acknowledgement is received. TCP only determines that connectivity has been restored after a (possibly long) retransmit timeout followed by the successful receipt of an ACK. Thus a link disconnection can result in a long delay in sending accompanied by a significant reduction in the sending rate [AP99].



Measurements: End-to-end performance under realistic topology and routing policies can be studied; [WMWGB06] suggests controlling routing events by injecting well-designed routing updates at known times to emulate link failures and repairs.

#### **16.5. Changes in Wireless Links: Mobility**

Definition: Network and node mobility, both wired and wireless, allows users to roam from one network to another seamlessly without losing service.

Determining factors: Mobility is a key attribute of wireless networks. Mobility can be determined by the presence of intersystem handovers, an intrinsic property of most wireless links [HS03].

Effect on Congestion Control Metrics: Mobility presents a major challenge to transport protocols through the packet losses and delay introduced by handovers. In addition to delay and losses, handovers can also cause a significant change in link bandwidth and latency. Host mobility increases packet delay and delay variation, and also degrades the throughput of TCP connections in wireless environments. Also, in the event of a handoff, slowly-responsive congestion control can require considerable time to adapt to changes. For example a flow under-utilizes a fast link after a handover from a slow link [HS03].

Measurements: Further references are needed to specify how mobility is actually measured [JEAS03].

#### **17. Using the Tools Presented in this Document**

[To be done.]

#### **18. Related Work**

[Cite "On the Effective Evaluation of TCP" by Allman and Falk.]

#### **19. Conclusions**

[To be done.]

#### **20. Security Considerations**

There are no security considerations in this document.



## **21. IANA Considerations**

There are no IANA considerations in this document.

## **22. Acknowledgements**

Thanks to Xiaoliang (David) Wei for feedback and contributions to this document. The sections on "Challenging Lower Layers" and "Network Changes Affecting Congestion" are contributions from Jasani Rohan with Julie Tarr, Tony Desimone, Christou Christos, and Vemulapalli Archana. The section on "Congestion Control Mechanisms for Traffic, along with Sender and Receiver Buffer Sizes" includes contributions from Sara Landstrom.

### Informative References

- [MAWI] M.W. Group, Mawi working group traffic archive, URL "<http://tracer.csl.sony.jp/mawi/>".
- [A00] M. Allman, A Web Server's View of the Transport Layer, Computer Communication Review, 30(5), October 200.
- [AAR03] Alhussein A. Abouzeid, Sumit Roy, Stochastic Modeling of TCP in Networks with Abrupt Delay Variations, Wireless Networks, V.9 N.5, 2003.
- [AGR00] M. Allman, J. Griner, and A. Richard. TCP Behavior in Networks with Dynamic Propagation Delay. Globecom 2000, November 2000.
- [AKM04] B. Appenzeller, I. Keslassy, and N. McKeown, Sizing Router Buffers, SIGCOMM 2004.
- [AKSJ03] J. Aikat, J. Kaur, F.D. Smith, and K. Jeffay, Variability in TCP Roundtrip Times, ACM SIGCOMM Internet Measurement Conference, Maimi, FL, October 2003, pp. 279-284.
- [ANP06] On Monitoring of End-to-End Packet Reordering over the Internet, B. Ye, A. P. Jayasumana, and N. M. Piratla, 2006.
- [AP99] M. Allman and V. Paxson. On Estimating End-to-end Network Path Properties. SIGCOMM, September 1999.
- [BH02] F. Baccelli and D. Hong, AIMD, Fairness and Fractal Scaling of TCP Traffic, Infocom 2002.
- [BA05] E. Blanton and M. Allman, "On the Impact to Bursting on TCP Performance", Passive and Active Measurement Workshop, March



2005.

- [BPS99] Bennett, J. C. R., Partridge, C. and Shectman, N., "Packet Reordering is Not Pathological Network Behavior," Trans. on Networking IEEE/ACM, Dec. 1999, pp.789-798.
- [CBC95] C. Cunha, A. Bestavros, and M. Crovella, "Characteristics of WWW Client-based Traces", BU Technical Report BUCS-95-010, 1995.
- [CR04] M. Chan and R. Ramjee, (2004) Improving TCP/IP Performance over Third Generation Wireless Networks. IEEE Infocom 2004.
- [Dummynet] L. Rizzo, Dummynet, URL  
"[http://info.iet.unipi.it/~luigi/ip\\_dummynet/](http://info.iet.unipi.it/~luigi/ip_dummynet/)".
- [F02] C. J. Fraleigh, Provisioning Internet Backbone Networks to Support Latency Sensitive Applications. PhD thesis, Stanford University, Department of Electrical Engineering, June 2002.
- [FJ92] S. Floyd and V. Jacobson, On Traffic Phase Effects in Packet-Switched Gateways, Internetworking: Research and Experience, V.3 N.3, September 1992, p.115-156.
- [FK02] S. Floyd and E. Kohler, Internet Research Needs Better Models, Hotnets-I, October 2002.
- [GF04] A. Gurtov and S. Floyd. Modeling Wireless Links for Transport Protocols. ACM CCR, 34(2):85-96, April 2004.
- [GK98] B. Gavish and J. Kalvenes. The Impact of Satellite Altitude on the Performance of LEOS based Communication Systems. Wireless Networks, 4(2):199--212, 1998.
- [GM04] L. Grieco and S. Mascolo, Performance Evaluation and Comparison of Westwood+, New Reno, and Vegas TCP Congestion Control, CCR, April 2004.
- [GPAR02] A. Gurtov, M. Passoja, O. Aalto, and M. Raitola. Multi-layer Protocol Tracing in a GPRS Network. In Proc. of the IEEE Vehicular Technology Conference (VTC02 Fall), Sept. 2002.
- [H95] Huitema, C., (1995) Routing in the Internet. Prentice Hall PTR, 1995.
- [HK99] T. Henderson and R. Katz, (1999) Transport Protocols for Internet-Compatible Satellite Networks. IEEE Journal on Selected Areas in Communications, Vol. 17, No. 2, pp. 345-359, February 1999.





- [HMTG01] C. Hollot, V. Misra, D. Towsley, and W. Gong, On Designing Improved Controllers for AQM Routers Supporting TCP Flows, IEEE Infocom, 2001.
- [HS03] R. Hsieh and A. Seneviratne. A Comparison of Mechanisms for Improving Mobile IP Handoff Latency for End-to-end TCP. MOBICOM, Sept. 2003.
- [IMD01] G. Iannaccone, M. May, and C. Diot, Aggregate Traffic Performance with Active Queue Management and Drop From Tail. SIGCOMM Comput. Commun. Rev., 31(3):4-13, 2001.
- [IMLGK03] H. Inamura, G. Montenegro, R. Ludwig, A. Gurtov and F. Khafizov. TCP over Second (2.5G) and Third (3G) Generation Wireless Networks. [RFC 3484](#), IETF, February 2003.
- [JD02] H. Jiang and C. Dovrolis, Passive Estimation of TCP Round-trip Times, Computer Communication Review, 32(3), July 2002.
- [JEAS03] A. Jardosh, E. Belding-Royer, K. Almeroth, and S. Suri. Towards Realistic Mobility Models for Mobile Ad Hoc Networks. MOBICOM, Sept. 2003.
- [LC01] J. Liu and Mark Crovella, Using Loss Pairs to Discover Network Properties, ACM SIGCOMM Internet Measurement Workshop, 2001.
- [LC05] X. Luo and R. K. C. Chang, Novel Approaches to End-to-end Packet Reordering Measurement, 2005.
- [LMJ96] Labovitz, C., Malan, G.R., and Jahanian, F., (1996) Internet Routing Instability. Proceedings of SIGCOMM 96.
- [MAF05] A. Medina, M. Allman, and A. Floyd. Measuring the Evolution of Transport Protocols in the Internet. Computer Communication Review, April 2005.
- [NISTNet] NIST Net, URL "<http://snad.ncsl.nist.gov/itg/nistnet/>".
- [NM01] J. Neale and A. Mohsen. Impact of CF-DAMA on TCP via Satellite Performance. Globecom, Nov. 2001.
- [PF01] J. Padhye and S. Floyd, Identifying the TCP Behavior of Web Servers, SIGCOMM 2001, August 2001.
- [P96] Paxson, V., (1996) End-to-end Routing Behavior in the Internet. Proceedings of SIGCOMM 96, pp. 25-38, August 1992.



- [P97] V. Paxson. End-to-end Routing Behavior in the Internet. IEEE/ACM Transactions on Networking, 5(5):60115, October 1997.
- [PJD04] R. Prasad, M. Jain, and C. Dovrolis, On the Effectiveness of Delay-Based Congestion Avoidance, PFLDnet 2004, February 2004.
- [QZK01] L. Qiu, Y. Zhang, and S. Keshav, Understanding the Performance of Many TCP Flows, Comput. Networks, 37(3-4):277-306, 2001.
- [RSB01] R. Riedi, S. Sarvotham, and R. Varaniuk, Connection-level Analysis and Modeling of Network Traffic, SIGCOMM Internet Measurement Workshop, 2001.
- [RTW05] G. Raina, D. Towsley, and D. Wischik, Control Theory for Buffer Sizing, CCR, July 2005.
- [LS06] D. Leith and R. Shorten, Impact of Drop Synchronisation on TCP Fairness in High Bandwidth-Delay Product Networks, Proc. Protocols for Fast Long Distance Networks, 2006.
- [TG] Traffic Generators for Internet Traffic Web Page, URL "<http://www.icir.org/models/trafficgenerators.html>".
- [UA01] U. of Auckland, Auckland-vi trace data, June 2001. URL "<http://wans.cs.waikato.ac.nz/wand/wits/auck/6/>".
- [UW02] UW-Madison, Network Performance Statistics, October 2002. URL "<http://wwwstats.net.wisc.edu/>".
- [W01] B. Walke. Mobile Radio Networks, Networking and Protocols (2. Ed.). Wiley & Sons, 2001.
- [WM05] D. Wischik and N. McKeown, Buffer sizes for Core Routers, CCR, July 2005. URL "<http://yuba.stanford.edu/~nickm/papers/BufferSizing.pdf>".
- [WMWGB06] F. Wang, Z. M. Mao, J. Wang, L. Gao and R. Bush. A Measurement Study on the Impact of Routing Events on End-to-End Internet Path Performance, SIGCOMM, 2006.
- [ZM04] X. Zhou and P. Van Mieghem, Reordering of IP Packets in Internet, PAM 2004.
- [ZSC91] L. Zhang, S. Shenker, and D.D. Clark, Observations and Dynamics of a Congestion Control Algorithm: the Effects of Two-way Traffic, SIGCOMM 1991.



[22]

[23]

[24]

#### Authors' Addresses

Sally Floyd  
ICSI Center for Internet Research  
1947 Center Street, Suite 600  
Berkeley, CA 94704  
USA  
Email: floyd@acm.org

Eddie Kohler  
4531C Boelter Hall  
UCLA Computer Science Department  
Los Angeles, CA 90095  
USA  
Email: kohler@cs.ucla.edu

#### Full Copyright Statement

Copyright (C) The IETF Trust (2008). This document is subject to the rights, licenses and restrictions contained in [BCP 78](#), and except as set forth therein, the authors retain all their rights.

This document and the information contained herein are provided on an "AS IS" basis and THE CONTRIBUTOR, THE ORGANIZATION HE/SHE REPRESENTS OR IS SPONSORED BY (IF ANY), THE INTERNET SOCIETY, THE IETF TRUST AND THE INTERNET ENGINEERING TASK FORCE DISCLAIM ALL WARRANTIES, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY THAT THE USE OF THE INFORMATION HEREIN WILL NOT INFRINGE ANY RIGHTS OR ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE.

#### Intellectual Property

The IETF takes no position regarding the validity or scope of any Intellectual Property Rights or other rights that might be claimed to pertain to the implementation or use of the technology described in this document or the extent to which any license under such rights might or might not be available; nor does it represent that it has made any independent effort to identify any such rights. Information on the procedures with respect to rights in RFC documents can be found in [BCP 78](#) and [BCP 79](#).



Copies of IPR disclosures made to the IETF Secretariat and any assurances of licenses to be made available, or the result of an attempt made to obtain a general license or permission for the use of such proprietary rights by implementers or users of this specification can be obtained from the IETF on-line IPR repository at <http://www.ietf.org/ipr>.

The IETF invites any interested party to bring to its attention any copyrights, patents or patent applications, or other proprietary rights that may cover technology that may be required to implement this standard. Please address the information to the IETF at [ietf-ipr@ietf.org](mailto:ietf-ipr@ietf.org).

