

The Performance of TCP/IP for Networks with High Bandwidth-Delay Products and Random Loss

T. V. Lakshman, *Member, IEEE*, and Upamanyu Madhow, *Senior Member, IEEE*

Abstract—This paper examines the performance of TCP/IP, the Internet data transport protocol, over wide-area networks (WANs) in which data traffic could coexist with real-time traffic such as voice and video. Specifically, we attempt to develop a basic understanding, using analysis and simulation, of the properties of TCP/IP in a regime where: 1) the bandwidth-delay product of the network is high compared to the buffering in the network and 2) packets may incur random loss (e.g., due to transient congestion caused by fluctuations in real-time traffic, or wireless links in the path of the connection). The following key results are obtained. First, random loss leads to significant throughput deterioration when the product of the loss probability and the square of the bandwidth-delay product is larger than one. Second, for multiple connections sharing a bottleneck link, TCP is grossly unfair toward connections with higher round-trip delays. This means that a simple first in first out (FIFO) queueing discipline might not suffice for data traffic in WANs. Finally, while the recent Reno version of TCP produces less bursty traffic than the original Tahoe version, it is less robust than the latter when successive losses are closely spaced. We conclude by indicating modifications that may be required both at the transport and network layers to provide good end-to-end performance over high-speed WANs.

Index Terms—Flow control, congestion control, error recovery, Internet, TCP/IP, transport protocols.

I. INTRODUCTION

MOST existing data transfer protocols have been designed for local-area network (LAN) applications in which buffer sizes far exceed the bandwidth-delay product.¹ This assumption may not hold for the wide-area networks (WANs) formed by the interconnection of LANs using high-speed backbone networks. In addition, in the Internet of the future, data traffic will share the network with voice and video traffic. In this paper, we examine the impact of these changes on the performance of the most popular data transfer protocol in current use, TCP/IP. This is essential not only for network provisioning in the short term (since the rapid growth of Web applications has caused TCP traffic to grow correspondingly)

Manuscript received June 20, 1995 revised February 25, 1997; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor D. Mitra. This work was supported in part by the U.S. Army Research Office under Grant DAAH04-95-1-0246.

T. V. Lakshman is with the High Speed Networks Research Dept., Bell Laboratories, Lucent Technologies, Holmdel, NJ 07733 USA (e-mail: lakshman@research.bell-labs.com).

U. Madhow is with the ECE Department and the Coordinated Science Laboratory, University of Illinois, Urbana, IL 61801 USA (e-mail: madhow@uiuc.edu).

Publisher Item Identifier S 1063-6692(97)04489-0.

¹The *bandwidth-delay product* is loosely defined to be the product of the round-trip delay for a data connection and the capacity of the bottleneck link in its path.

but also for determining how TCP needs to be modified in the longer term.

We study two versions of TCP: one is the popular Tahoe version developed by Jacobson in 1988 [11] (henceforth called *TCP-tahoe*); the other is the Reno version, which includes the fast retransmit option together with a method for reducing the incidence of slow start, suggested by Jacobson in 1990 [12] (we will refer to this as *TCP-reno*). We attempt to develop a basic understanding of these schemes by considering one-way traffic over a single bottleneck link with FIFO transmission. For LANs, the round-trip delay of the connection is small, so that the bandwidth-delay product could be much smaller than the buffering on the bottleneck link. We are more interested, however, in WANs with large round-trip delays, so that the buffering on the bottleneck link is typically of the same order of magnitude as, or smaller than, the bandwidth-delay product (this is what we mean by *high bandwidth-delay products* throughout this paper). The bottleneck link may be shared by several TCP connections. In addition, we also assume that each packet may be lost randomly even after obtaining service at the bottleneck link.

Random loss is a simple model for a scenario of particular interest in the context of networks with multimedia traffic, where transient fluctuations in real time traffic may cause irregularly spaced losses for data traffic. This would occur, for instance, for both the UBR and ABR service classes [1] in ATM networks. The only difference is that for ATM ABR, each connection would have a time-varying available rate determined by feedback from the switches, so that most random losses would occur at the interface of the source to the network, since that is where the available rate would be enforced. In addition to serving as a model for transient congestion, we note that random loss on the Internet has been reported [3], where it is conjectured to occur due to a variety of reasons, including intermittent faults in hardware elements such as Ethernet/FDDI adapters, and incorrect handling of arriving packets by routers. Finally, with the anticipated emergence of mobile computing over heterogeneous networks with both wireless and wireline links, losses and time variations due to wireless links in the path of the connection can also be accommodated via a random loss model. Since our purpose is to obtain a fundamental understanding of TCP, none of the preceding situations are explicitly considered in this paper. However, as discussed in Section VI, the results here should provide a basis for further work on developing network level design guidelines for supporting TCP.

One of the simplifications of the model used for our analysis is that two-way traffic (and the accompanying *ack compression*

[27]) is not considered. Feedback systems are notoriously difficult to analyze, so that even our simple model is not amenable to exact analysis. However, not only does our approximate analysis match simulation results for the idealized system model, but it also provides a close match to results for a detailed simulation that includes two-way traffic for multiple TCP-Reno connections over an ATM network (described in Section V).

We obtain the following key results. Discussion of the implications of these results for system design is postponed to Section VI.

- 1) While TCP-reno produces less bursty traffic than TCP-tahoe, it is much less robust toward “phase effects.” The latter term refers to unpredictability in performance resulting from very small differences in the relative timings of packet arrivals for different connections sharing a link.
- 2) Both versions of TCP appear to have significant drawbacks as a means of providing data services over multimedia networks, because random loss resulting from fluctuations in real-time traffic can lead to significant throughput deterioration in the high bandwidth-delay product regime. Roughly speaking, the performance is degraded when the product of the loss probability and the *square* of the bandwidth-delay product is large (e.g., ten or more).
- 3) For high bandwidth-delay products, TCP is grossly unfair toward connections with higher propagation delays: for multiple connections sharing a bottleneck link, the throughput of a connection is inversely proportional to (a power of) its propagation delay.

It is worth clarifying that random loss causes performance deterioration in TCP because it does not allow the TCP window to reach high enough levels to permit good link utilization. On the other hand, when the TCP window is already large and is causing congestion, random early drops of packets when the link buffer gets too full can actually enhance performance and alleviate phase effects [10].

Early simulation studies of TCP-tahoe include [24], [26], [27]. Our model is similar to that used in [24], but the key differences between our paper and previous studies are that: 1) the ratio of bandwidth-delay product to buffer size is much higher in our study and 2) the effect of random loss due to transient congestion (or other sources) is included. Thus, some of the undesirable features of TCP-tahoe which arise specifically for networks with high bandwidth-delay products (such as excessive buffering requirements and vulnerability to random loss) were not noticed in earlier studies. Furthermore, in contrast to previous studies, we place more emphasis on detailed analytical insight on the effects of various parameters on performance.

The bias of TCP-tahoe against connections with large round-trip delays and against connections traversing a large number of congested gateways has been noticed in other studies of TCP-tahoe [8], [9], [26]. A heuristic analysis in [8] shows that, for multiple connections sharing a bottleneck link, the throughput of a connection is inversely proportional to its

round-trip time. While we consider a similar system in Section V, our analysis is more detailed, taking explicit account of the buffer size and the bandwidth-delay product. Oscillatory behavior and unfairness toward connections with larger propagation delays have also been noticed in a previous analytical study of feedback-based congestion control [2]. Other analyses of flow control schemes include [20], [22], [23], but these references do not address the specific concerns raised here in any detail.

The system model is described in Section II. Analytical and simulation results for the evolution of a single connection in the absence of random loss are given in Section III. Section IV considers the effect of random loss. Section V contains results for multiple connections with and without random loss. We give our conclusions in Section VI.

II. SYSTEM MODEL

We consider infinite data sources which always have packets to send, so that the units of data are maximum sized packets (in general, packet sizes in TCP may be variable). We consider a single bottleneck link with capacity μ packets per second and a FIFO buffer of size B packets. Any packet arriving when the buffer is full is lost (random loss may cause additional losses). The number of connections, or sources, sharing the link is assumed to be constant. For each connection, all delays except for service time and queueing at the bottleneck link are lumped into a single “propagation delay,” which includes: 1) the time between the release of a packet from the source and its arrival into the link buffer; 2) the time between the transmission of the packet on the bottleneck link and its arrival at its destination; and 3) the time between the arrival of the packet at the destination and the arrival of the corresponding acknowledgment at the source. The propagation delay for a packet from the i th connection is denoted by τ_i .

The τ_i are taken to be deterministic, which implicitly assumes that deterministic propagation and processing delays are more significant than random queueing delays at all nodes and links other than the bottleneck link. Although such an assumption is overly simplistic even for a relatively simple system with two-way traffic [27], it suffices for our present purpose of arriving at a basic understanding of the interaction between different connections sharing a link.

Each connection is assumed to use a window flow control protocol. At time t , the window size for connection i is denoted by $W_i(t)$, and is equal to the maximum allowed number of unacknowledged packets (retransmissions are not counted). It is assumed that each connection uses its allowable window to the fullest extent, i.e., that at time t , there are indeed $W_i(t)$ unacknowledged packets for connection i . The window varies dynamically in response to acknowledgment and detection of packet loss. Upon receiving a packet, the destination is assumed to send an acknowledgment back immediately. These acknowledgments are *cumulative* and indicate the next byte expected by the receiver.

In the original version of TCP-tahoe, packet loss is detected by maintaining a timer based on an estimate of the round-trip time. When a packet is sent, a timeout value is computed using the current round-trip time estimate and the timer is

started. Expiry of this timer is taken to signal packet loss. For each retransmission following a timer expiry, the timer value used is twice the previous timer value. Estimates of the round-trip time are obtained by measuring the round-trip time upon receipt of unambiguous acknowledgment (i.e., ignoring acknowledgment for retransmitted segments) and computing a weighted average of the old and new estimates. Refer to [15], [25] for a detailed description of round-trip time estimation. We will refer to a timer based on this estimate as a *fine-grained timer*, in order to distinguish it from the *coarse-grained timers* used in practice, which are typically multiples of 500 ms. In order to prevent a needlessly lengthy stoppage of transmission upon expiry of a coarse-grained timer, most current versions of both TCP-tahoe and TCP-reno incorporate a *fast retransmit* option: if the number of *duplicate acknowledgments* (i.e., multiple acknowledgment with the same “next expected” packet number n) exceeds a threshold, packet n is assumed to be lost. In this paper, we implement fine-grained timers in our simulations, in order to study the dynamic evolution of TCP (and to highlight possible shortcomings) in the most ideal setting. The original version of TCP-tahoe, without fast retransmit, is implemented. However, in simulation results not reported here, we have checked that coarse-grained timers with fast retransmit give virtually identical performance in most cases of interest for TCP-tahoe (unless almost all packets in a window are lost, fast retransmit detects loss very effectively). For TCP-reno, we implement fast retransmit with a fine-grained timer in our simulations. Because TCP-reno has a less robust congestion control mechanism, we have found in later work that the use of a coarse-grained timer does impact its performance even with fast retransmit, unlike for TCP-tahoe. Since either fine-grained timers or the fast retransmit option provide almost perfect loss detection, it is assumed in our analysis that packet losses are detected perfectly.

A simplified description of TCP-tahoe [11] and TCP-reno [12] follows.

A. Description of TCP-tahoe

The algorithm followed by each connection has two parameters, current window size W and a threshold W_t , which are updated as follows.

- 1) After every acknowledgment of a new packet:
if $W < W_t$, set $W = W + 1$; **Slow Start Phase**
else set $W = W + 1/[W]$. **Congestion Avoidance Phase**
($[x]$ denotes the integer part of x).
- 2) After a packet loss is detected:
set $W_t = W/2$;
set $W = 1$.

The algorithm typically evolves as follows (although, as described in the next section, the evolution is somewhat different for relatively small buffer size): when packet loss is detected, the window is reduced to one. In the slow start phase that follows, the window grows rapidly for every successfully acknowledged packet until it reaches half of the window size

at the last packet loss. The algorithm then switches to the congestion avoidance phase, probing for extra bandwidth by incrementing the window size by one for every window’s worth of acknowledged packets. This growth continues until another packet loss is detected, at which point another cycle begins. We use the term *cycle* to mean TCP evolution starting from the end of one congestion avoidance phase to the end of the next. In Section III, it turns out that, for our simple model, TCP evolution is periodic if there is no random loss, so that successive cycles are identical. In Section IV, on the other hand, where we consider random loss, the duration of, and window evolution within, different cycles is random.

B. Description of TCP-reno

After the number of duplicate acknowledgments exceeds a threshold (typically three), TCP-reno retransmits the packet. However, instead of cutting the window back to one, it only reduces it by a factor of two. Further, in order to prevent a burst of packets from being transmitted when the retransmission is finally acknowledged, it temporarily permits new packets to be transmitted with each repeated acknowledgment until the “next expected” number in the acknowledgment advances. While these subtleties are essential to the working of the algorithm (see [12] for details) and are implemented in our simulations, the following simplified description is adequate for conveying an understanding of the algorithm’s behavior.

- 1) After every nonrepeated acknowledgment, the algorithm works as before:
if $W < W_t$, set $W = W + 1$; **Slow Start Phase**
else set $W = 1 + 1/[W]$. **Congestion Avoidance Phase.**
- 2) When the duplicate acknowledgment exceeds a threshold,
retransmit “next expected” packet;
set $W_t = W/2$, then set $W = W_t$ (i.e., halve the window);
resume congestion avoidance using new window once retransmission is acknowledged.
- 3) Upon timer expiry, the algorithm goes into slow start as before:
set $W_t = W/2$;
set $W = 1$.

In this case, after an initial slow start transient, the typical cyclical evolution does not involve slow start, since the window size is halved upon loss detection. Each cycle begins when a loss is detected via duplicate acknowledgment. Assuming that loss occurs at window size W_{\max} , the window size at the beginning of each cycle is $W_{\max}/2$. The algorithm resumes probing for excess bandwidth in congestion avoidance mode until the window size reaches W_{\max} again, at which point a loss occurs and a new cycle with window size $W_{\max}/2$ begins. We will show that the throughput attained by this scheme is higher than that of TCP-tahoe, especially when the buffer size is small compared to the bandwidth-delay product. However, this algorithm is almost as vulnerable to random loss.

For the remainder of this paper, we will use W_{\max} as a generic notation for the window size at which congestion

avoidance ends. The value of W_{\max} could therefore change from cycle to cycle if loss occurs randomly, or could be the same for all cycles if loss occurs periodically. It is worth relating our notation to that usually used in TCP code (see [24], for instance): W is usually referred to as the congestion window $cwnd$, and W_t is denoted by $ssthresh$. The actual window for flow control purposes is taken to be the minimum of $cwnd$ and $maxwnd$, where the latter is set by the receiver. For the purpose of this paper, the window size is assumed to be dictated by the capacity and buffering of the bottleneck link (i.e., $cwnd \leq maxwnd$), so the actual window size equals the congestion window. Note that some form of window scaling (i.e., increasing the window size in bytes while using the same sequence number space, by scaling up the size of the data segment referred to by a given number) may be required to achieve this for large bandwidth-delay products [14].

III. EVOLUTION WITHOUT RANDOM LOSS

We consider the evolution of a single connection and derive expressions for its long-term throughput. Define the normalized buffer size $\beta = B/(\mu\tau + 1) = B/\mu T$, where τ denotes the propagation delay for each packet of the connection and $T = \tau + 1/\mu$ denotes the propagation delay plus the service time. Since we are concerned with large bandwidth-delay products, we restrict attention to $\beta \leq 1$ in this section. In contrast, simulations in earlier work [24] consider $\beta \gg 1$, for which the average throughput is close to the capacity of the bottleneck link. For brevity, expressions for the latter case are omitted.

The maximum window size that can be accommodated in steady state in the bit pipe is

$$W_{\text{pipe}} = \mu T + B = \mu\tau + B + 1. \quad (1)$$

In this case, the buffer is always fully occupied and there are μT packets in flight. The cyclical evolution of TCP-tahoe consists of a slow start phase starting with $W = 1$ and continuing until the window size reaches $W_t = W_{\text{pipe}}/2$, followed by congestion avoidance until $W = W_{\text{pipe}}$. The next increase in window size leads to buffer overflow, at which point the window is reset to one and a new cycle starts. We show that if the relative buffer size β is not large enough, buffer overflow may occur even in the slow start phase, and the cyclical evolution is somewhat different from the preceding description. For TCP-reno, if the scheme functions as designed, slow start is eliminated from the cyclical evolution. In each cycle, the algorithm starts from $W = W_t = W_{\text{pipe}}/2$, does congestion avoidance until $W = W_{\text{pipe}}$, and drops back to $W = W_t = W_{\text{pipe}}/2$ after a packet loss due to buffer overflow is detected via duplicate acknowledgment.

In each case, if the number of packets successfully transmitted during a cycle is N_c and the duration of a cycle is T_c , then the periodic evolution implies that the average throughput is given by $\bar{\lambda} = N_c/T_c$. In the following, we describe this evolution more carefully, and compute these quantities in sufficient detail to produce an excellent match with simulations (see Table II).

TABLE I
EVOLUTION DURING SLOW START PHASE

Time	Packet Acked	Window Size	Packet(s) Released	Queue Length
0		1	1	1
T	1	2	2, 3	2
$2T$	2	3	4, 5	2
$2T + 1/\mu$	3	4	6, 7	$2 - 1 + 2 = 3$
$3T$	4	5	8, 9	2
$3T + 1/\mu$	5	6	10, 11	$2 - 1 + 2 = 3$
$3T + 2/\mu$	6	7	12, 13	$2 - 1 + 2 = 3$
				$1 + 2 = 4$
$3T + 3/\mu$	7	8	14, 15	$2 - 1 + 2 = 3$
				$1 + 2 - 1 + 2 = 5$
$4T$	8	9	16, 17	2
.				
.				

A. Slow Start Phase

The slow start phase must be considered in some detail to understand the advantage of TCP-reno over TCP-tahoe. Starting from $W = 1$ with slow start threshold W_t , the window size is increased by one for every acknowledgment in this phase, so that two packets are released into the buffer for every acknowledgment. Table I shows the evolution of the window size and the queue length in this phase. For every acknowledgment, we indicate the number of the packet which was acknowledged (for convenience, we number the packets in increments of one rather than in increments equal to the number of bytes per packet).

The evolution in Table I is best described by considering *mini-cycles* of duration equal to the round-trip time T , where the i th mini-cycle refers to the time interval $[iT, (i+1)T)$ (the mini-cycles are separated by lines in the table). The acknowledgment for a packet released in mini-cycle i arrives in mini-cycle $(i+1)$, and increases the window size by one. This leads to a doubling of the window in each mini-cycle. Further, acknowledgment for consecutive packets served in mini-cycle i arrive spaced by the service time $1/\mu$ during mini-cycle $(i+1)$, and two packets are released for each arriving acknowledgment, leading to a buildup of queue size. The preceding evolution assumes implicitly that the normalized buffer size $\beta < 1$, so that the window size during the slow start phase is smaller than μT and the queue empties out by the end of each mini-cycle. Denoting the window size at time t by $W(t)$, we obtain that, during the $(n+1)$ th mini-cycle,

$$W(nT + m/\mu) = 2^n + m + 1, \quad 0 \leq m \leq 2^n - 1 \quad (2)$$

where we have assumed that $(2^n - 1)/\mu < T$. Similarly, letting $Q(t)$ denote the queue length at time t , the queue build-up during the n th mini-cycle is given by

$$Q(nT + m/\mu) = m + 2, \quad 0 \leq m \leq 2^n - 1. \quad (3)$$

The maximum queue length during the $(n+1)$ th mini-cycle is therefore $2^n + 1$, which is approximately half the maximum window size $W[nT + (2^n - 1)/\mu] = 2^{n+1}$ during that mini-cycle. For a buffer size B , we can use (2) and (3) to determine the window size at which, the queue length exceeds the buffer

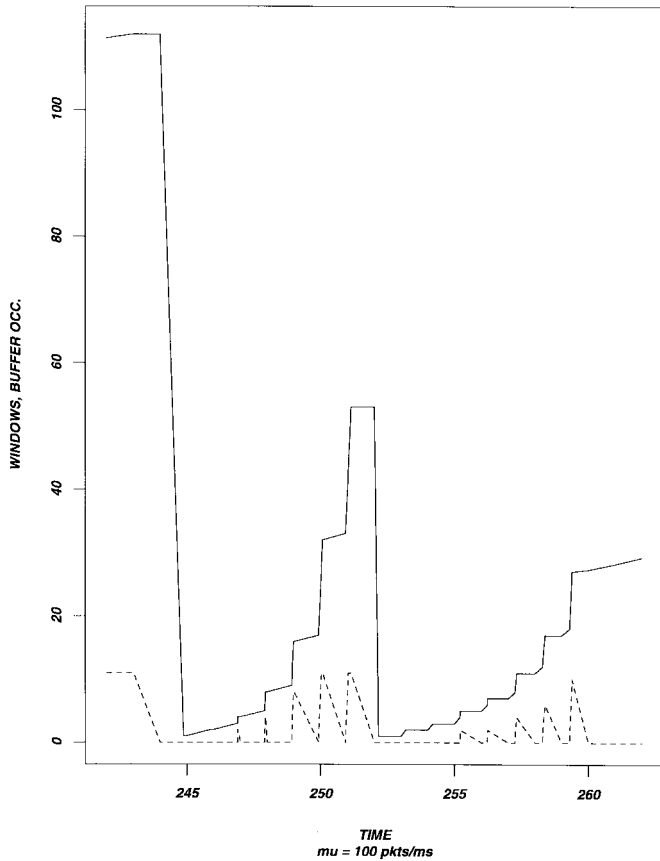


Fig. 1. Window and buffer evolution for a single connection: Two slow starts. Prop. delay = 1 ms; $b = 0.1$.

size as follows. Define the integer $n_b = \lceil \log_2(B - 1) \rceil$, so that $2^{n_b-1} + 1 < B \leq 2^{n_b} + 1$. From (3), buffer overflow will occur in the n_b th mini-cycle (the largest queue length in the previous cycle is $2^{n_b-1} + 1$, which is smaller than B), with $m + 2 = B + 1$. From (2), the window size W_b at which this happens is $2^{n_b} + m + 1$, so that that

$$W_b = 2^{n_b} + B. \quad (4)$$

Buffer overflow during a slow start phase with threshold W_t thus occurs only if

$$W_b \leq W_t. \quad (5)$$

A more explicit condition for buffer overflow can be derived as follows. Assuming that the packet loss causing the slow start phase occurred when the window size exceeds the value $W_{\text{pipe}} = \mu T + B$, the slow start threshold equals $W_t = W_{\text{pipe}}/2 = (1 + \beta)\mu T/2$. Since $W_b \approx 2B = 2\beta\mu T$, the condition for buffer overflow (5) is approximately equivalent to $\beta \leq 1/3$.

Fig. 1 shows the simulated congestion window and buffer occupancy evolution for a single connection using TCP-tahoe with $\tau = 1$, $\mu = 100$, and $\beta = 0.1$ (i.e., $B = 10$). The congestion window size is shown by the solid line and the buffer occupancy by the dotted line. As expected, the window grows to $W_{\text{pipe}} = \mu\tau + B + 1 = 111$ and the next increase in the window causes a packet to be dropped. Detection of this loss (upon expiry of the associated timer) causes the window to be reduced to one and initiation of the slow start phase. The

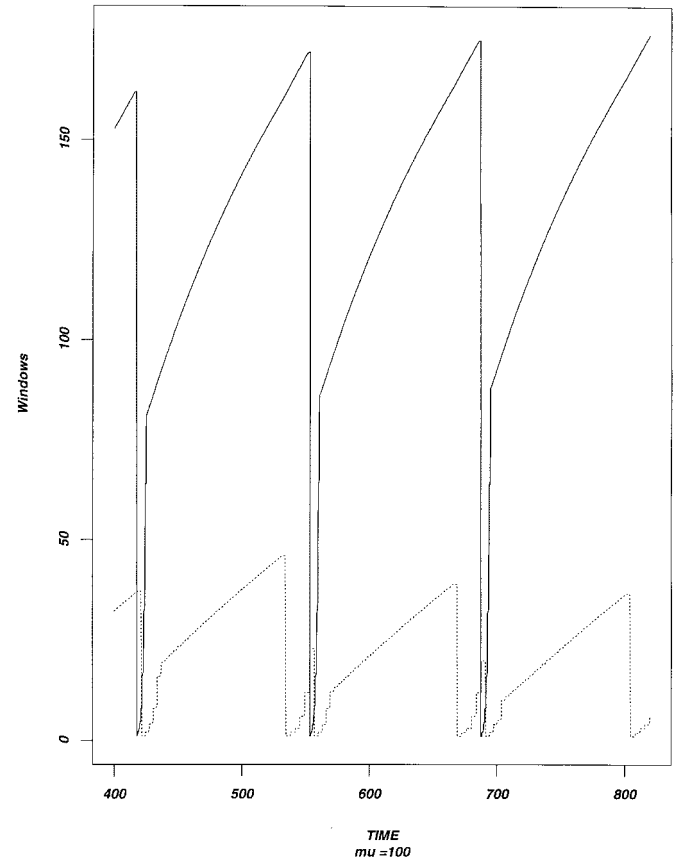


Fig. 2. Window and buffer evolution for a single connection: one slow start. $\tau_1 = 1, \tau_2 = 3, b = 0.8$.

figure clearly shows the rapid growth in window size during the slow start phase. However, since $\beta < 1/3$, buffer overflow occurs when $W = W_b = 26$, and is detected by the time the window size reaches $W = 2W_b - 2 = 50$ (see the discussion later in this section). A second slow start phase is initiated at this point with threshold 25. This window size is reached without further loss, at which point slower window growth due to congestion avoidance commences. This lasts until the window exceeds $W_{\text{pipe}} = 111$, after which a new cycle begins.

When β is greater than $1/3$ the window evolution for TCP-tahoe is different. This is illustrated in Fig. 2, which shows the evolution of window sizes and buffer occupancy for $\beta = 0.8$. Here, packet loss is seen to occur when the window is of size $W_{\text{pipe}} = 181$. As before, detection of this loss causes the window size to be reduced to one and initiates slow start. However, there is no packet loss in the slow start phase, which terminates when the window reaches 91. In the congestion avoidance phase that follows, the window grows linearly and then more slowly (as explained in the next subsection) until the window exceeds W_{pipe} . This results in a packet loss causing the cycle to repeat. The absence of the double slow start results in much higher throughput, since the initial window size for the congestion avoidance phase (which accounts for most of the packets transmitted) is higher.

We now compute the duration and number of packets transmitted during the slow start phase(s) in a given cycle for TCP-tahoe. Even though many subtleties in timing are glossed over, the computations are accurate enough to re-

sult in excellent agreement with the throughput obtained by simulations.

Case 1 $W_b > W_t$: There is only one slow start phase per cycle, which ends when the window size reaches $W_t = W_{\text{pipe}}/2$. We use a simplified version of (2), $W(t) = 2^{t/T}$ to approximate the duration of this phase as $t_{ss} = T \log_2 W_t$. The number of packets transmitted in this phase is approximated by $n_{ss} = W_t$. (Since the window size grows by one for every acknowledgment during slow start, starting from an initial value of one, the number of packets transmitted successfully during slow start is well approximated by the window size at the end of the slow start period.)

Case 2 $W_b \leq W_t$: There are two slow start phases in a cycle in this case. Let t_{ss1} denote the duration of the first slow start phase and let n_{ss1} denote the number of packets successfully transmitted during that phase. Let t_{ss2} , n_{ss2} denote the analogous quantities for the second slow start phase. Computation of average throughput requires the computation of these quantities.

In the first slow start phase (with threshold $W_t = W_{\text{pipe}}/2$), a buffer overflow occurs when the window size reaches W_b . The duration of this phase is the time taken to reach W_b , approximated by $T \log_2 W_b$, together with the time taken to detect the loss, which is taken to be one round-trip time T , so that $t_{ss1} = T(\log_2 W_b + 1)$. The number of packets transmitted in this phase is, reasoning as before, taken to be $n_{ss1} = W_b$.

Since the buffer overflow in the first slow start phase is not detected till the mini-cycle after which it happens, the window size at which the overflow is detected can be shown, by means of a more careful analysis, to be approximately $W^* \approx \min\{2W_b - 2, W_t\}$. This implies that the threshold for the second slow start phase is

$$\begin{aligned} \tilde{W}_t &= \frac{W^*}{2} \\ &\approx \min\left\{W_b - 1, \frac{W_t}{2}\right\}. \end{aligned} \quad (6)$$

The duration and number of packets transmitted during the second slow start phase is now obtained as in Case 1 to be $t_{ss2} = T \log_2 \tilde{W}_t$ and $n_{ss2} = \tilde{W}_t$, respectively.

The total time spent in slow start and the number of packets transmitted during slow start are then given by $t_{ss} = t_{ss1} + t_{ss2}$ and $n_{ss} = n_{ss1} + n_{ss2}$, respectively.

B. Congestion Avoidance Phase

We can unify the analysis of this phase for TCP-tahoe and TCP-reno by assuming that the congestion avoidance phase starts from an arbitrary window size W_0 and terminates when the window size reaches W_{max} . In all cases considered in this section, $W_{\text{max}} = W_{\text{pipe}}$. For TCP-tahoe, W_0 equals the slow start threshold for the slow start phase immediately preceding the congestion avoidance phase, and is given by

$$W_0(\text{tahoe}) = \begin{cases} \frac{W_{\text{max}}}{2}, & \frac{W_{\text{max}}}{2} < W_b \\ \min\left[W_b - 1, \frac{W_{\text{max}}}{4}\right], & \frac{W_{\text{max}}}{2} \geq W_b. \end{cases} \quad (7)$$

Remark: For $W_{\text{max}}/2 \geq W_b$, there are two slow start phases, so that the window size at the beginning of congestion avoidance is the slow start threshold for the second slow start phase, which is given in (6).

For TCP-reno, we have

$$W_0(\text{reno}) = \frac{W_{\text{max}}}{2} \quad (8)$$

since the window size is halved after losing a packet.

In contrast to the slow start phase, the congestion avoidance phase, by virtue of its slower window growth, is well-modeled by a continuous-time approximation for window evolution. Such approximations have been used previously in [24] to explain simulation results. Let dW/dt denote the rate of window growth with time, dW/da the rate of window growth with arriving acknowledgments, and da/dt the rate at which the acknowledgments are arriving. Then, during the congestion avoidance phase,

$$dW/da = 1/W. \quad (9)$$

The acknowledgments arrive back at a rate equal to the instantaneous throughput λ , so that

$$\frac{da}{dt} = \lambda = \min\left\{\frac{W}{T}, \mu\right\}. \quad (10)$$

Combining (9) and (10), we obtain that

$$\frac{dW}{dt} = \begin{cases} \frac{1}{T}, & W \leq \mu T \\ \frac{\mu}{W}, & W \geq \mu T. \end{cases} \quad (11)$$

Thus, for $W \leq \mu T$, the window W grows as t/T . The duration of this period of growth is therefore given by

$$t_A = T(\mu T - W_0) \quad (12)$$

since the initial window size was W_0 (for $\beta < 1$, $W_0 \leq W_{\text{max}}/2$ is always less than μT). The number of packets transmitted during this time is given by

$$n_A = \int_0^{t_A} \frac{W(t+t_{ss})}{T dt} = \int_0^{t_A} \frac{W_0 + \frac{t}{T}}{T dt} = \frac{W_0 t_A}{T} + \frac{t_A^2}{2T}. \quad (13)$$

When $W > \mu T$ ($t \geq t_A$), we obtain from (11) that $W^2(t) = 2\mu(t - t_A) + (\mu T)^2$. This growth period, and the cycle, terminates with buffer overflow when the window size exceeds W_{max} , and its duration is given by

$$t_B = \frac{W_{\text{max}}^2 - (\mu T)^2}{2\mu} \quad (14)$$

($t_B = 0$ if $W_{\text{max}} < \mu T$, although this does not occur for $W_{\text{max}} = W_{\text{pipe}} = \mu T + B$). Since the bottleneck link is being fully utilized during this period, the number of packets transmitted is given by

$$n_B = \mu t_B. \quad (15)$$

TABLE II
LINK UTILIZATION AS A FUNCTION OF
NORMALIZED BUFFER SIZE ($\mu = 100, \tau = 1$)

β	Link Utilization (Analysis/Simulation)	
	Tahoe	Reno
.1	.604/.604	.818/.818
.2	.660/.664	.871/.870
.31	.708/.718	.915/.911
.32	.856/.858	.919/.916
.8	.953/.954	.996/.994

C. Throughput Computation and Numerical Results

Due to the periodic evolution, the long-run average throughputs for both TCP-tahoe and TCP-reno are equal to the average throughputs in a cycle, and are given by

$$\bar{\lambda} = \frac{n_{ss} + n_A + n_B}{t_{ss} + t_A + t_B}, \quad \text{TCP-tahoe} \quad (16)$$

$$\bar{\lambda} = \frac{n_A + n_B}{t_A + t_B}, \quad \text{TCP-reno} \quad (17)$$

where the preceding quantities are as computed in Sections III-A and III-B.

Table II gives the link utilizations as a function of the normalized buffer size β for both TCP-tahoe and TCP-reno. The results obtained using (16) and (17) are compared with those obtained using simulation for an example with high bandwidth-delay product. The match is within 2%. For TCP-tahoe, a clear thresholding effect is seen at $\beta = 0.31$ (recall that the analysis predicted a thresholding effect around $\beta \approx 1/3$). The utilization for TCP-reno is uniformly higher for all values of β , and there is no thresholding effect for small β . The difference in utilization is small for large β , since the congestion avoidance phase for the two schemes is identical, and the duration of the slow start phase is small compared to the duration of the cycle.

IV. EVOLUTION WITH RANDOM LOSS

We assume here that any given packet may be lost with probability q , and that these random losses are independent. As in the previous section, we consider a single connection, and show that, for both TCP-tahoe and TCP-reno, the throughput is strongly dependent on $q(\mu T)^2$, and deteriorates sharply compared to the lossless throughput when this quantity becomes large. Roughly speaking, the throughput degradation occurs because packet losses relatively early in a cycle result in small initial values for the congestion avoidance phase, in which the bulk of the packets are transmitted. This results in small window sizes (determined by random losses rather than congestion), and therefore low link utilizations, during the cycle.

In principle, it is possible to exactly compute the throughput based on a Markov chain analysis, and we sketch this method in the following. However, insight into when random loss causes throughput deterioration is better obtained by means of an approximate analysis, which we pursue in more detail.

In the absence of random loss, the evolution of a cycle in TCP-tahoe is completely determined by the slow start threshold W_t (which is half the window size at the end of the previous cycle). Similarly, the evolution of a cycle in TCP-

reno is determined by the window size W_0 at the beginning of the cycle (again, this is half the window size at the end of the previous cycle). Since a single parameter w ($w = W_t$ for TCP-tahoe and $w = W_0$ for TCP-reno) determines the cyclical evolution, the following functions are well-defined for either scheme (although they may be hard to compute explicitly):

$W(n, w)$ window size after n packets are successfully transmitted;

$T(n, w)$ time taken for n packets to be successfully transmitted;

$N_{\max}(w)$ number of packets successfully transmitted before the cycle ends with buffer overflow.

If we now introduce random loss, the cycle may terminate due to a random loss after $N \leq N_{\max}$ packets have been successfully transmitted. According to our model of random loss, the probability distribution of N is specified by

$$P[N = n] = \begin{cases} (1-q)^n q, & n < N_{\max} \\ (1-q)^{N_{\max}}, & n = N_{\max}. \end{cases} \quad (18)$$

The window size when the cycle terminates is $W(N, w)$ and the duration of the cycle is $T(N, w)$. The cyclical evolution can now be completely specified as follows. For the m th cycle, let w_m , N_m , and T_m denote, respectively, the window size at the beginning of the congestion avoidance phase, the number of successful transmissions, and the duration. The threshold w_0 for the zeroth cycle is assumed to be given. According to our model for random loss, the random variables N_m are independent and identically distributed according to (18). The evolution of w_m and T_m is then specified as follows:

$$w_{m+1} = W(N_m, w_m) \quad (19)$$

$$T_{m+1} = T(N_m, w_m). \quad (20)$$

Equations (18) and (19) specify the transition probabilities for the time-homogeneous Markov chain formed by the $\{w_m\}$. This can be solved to obtain the stationary distribution for the w_m . The long-run throughput is then given by $\bar{\lambda} = E[N_m]/E[T_m]$.

Specifying the deterministic functions $W(n, w)$ and $T(n, w)$ in detail and solving for the stationary distribution is probably more time-consuming than simulation, and is not likely to yield any additional insight. In order to develop an intuitive understanding of how the random loss probability affects the throughput, therefore, we use an approximation for throughput based on the assumption that every cycle begins with a single ‘‘average’’ value $w = w^*$. This yields

$$\bar{\lambda} \approx \sum_n P[N = n] \frac{n}{T(n, w^*)}. \quad (21)$$

The parameter w^* is taken to be the minimum of its value without random loss, $W_{\text{pipe}}/2 = (\mu T + B)/2$ (this value is the same for TCP-tahoe and TCP-reno), and a value w_l based only on random loss, and computed in a simple-minded fashion as follows. For a loss probability q , roughly $1/q$ packets are successful in a cycle before a loss occurs. The parameter w is chosen to satisfy the following ‘‘fixed point’’ relationship: given that the parameter for the current cycle is w_l , and that exactly $1/q$ packets are successful, choose w_l such that the

window size at the time of loss is $2w_l$, so that the parameter for the next cycle is also w_l . Thus, w_l satisfies

$$\frac{W\left(\frac{1}{q}, w_l\right)}{2} = w_l. \quad (22)$$

Since the preceding approximation is fairly drastic, there is no point in solving (22) exactly, so that we feel free to resort to further approximations. Consider TCP-reno first. Starting from $W_0 = w_l$, suppose that the packet loss occurs at time t after the cycle begins. Assuming that the window size does not reach μT (i.e., that random loss has a significant effect in limiting window size), the linear growth in the congestion avoidance phase [see (18)] implies that $2w_l = w_l + t/T$, so that

$$t = Tw_l. \quad (23)$$

The number of packets transmitted in time t in the congestion avoidance phase is [see (20)] $[w_l t + t^2/(2T)]/T$, so that

$$\frac{w_l t + \frac{t^2}{2T}}{T} = \frac{1}{q}. \quad (24)$$

Combining (23) and (24), we obtain upon simplification that

$$w_l = \sqrt{2/3q}. \quad (25)$$

Starting with a slow start threshold $W_t = w_l$, the evolution for TCP-tahoe is the same as for TCP-reno with $W_0 = w_l$ if one ignores the slow start phase, which is of relatively short duration. The value for w_l obtained by solving (22) for TCP-tahoe is therefore also taken to be given by (25). The way this value is used in subsequent computations is somewhat different, however.

For both TCP-tahoe and TCP-reno, the window size at the beginning of the congestion avoidance phase strongly influences the throughput in the cycle, since the congestion avoidance phase accounts for the bulk of the packets transmitted in the cycle (the slow start phase is relatively short). Thus, in order for the throughput with random loss to be comparable with that without loss, the threshold

$$w^* = \min \left\{ \frac{\mu T + B}{2}, w_l \right\} \quad (26)$$

must be comparable (of the same order as) $(\mu T + B)/2$, the threshold in the absence of random loss. From (25), this holds if $q(\mu T + B)^2 = q(\mu T)^2(1 + \beta)^2$ is comparable to or smaller than $8/3$. We therefore expect a significant deterioration of throughput for both TCP-tahoe and TCP-reno when $q(\mu T)^2$ is large (say 10 or more). This is the case in Fig. 3 which shows the evolution of window sizes and buffer occupancies with $q = 0.001$ for the single connection whose lossless evolution was shown in Fig. 2. As expected, the slow start thresholds are much smaller than the thresholds for the lossless case (which is 91) and this results in severe throughput deterioration. In the time period shown in the figure, the buffer never overflows and all the packet losses are due to random losses.

Having specified w^* , it remains to compute the approximation (21). Since computation of the deterministic function $T(n, w^*)$ and summation over n is quite tedious, we resort

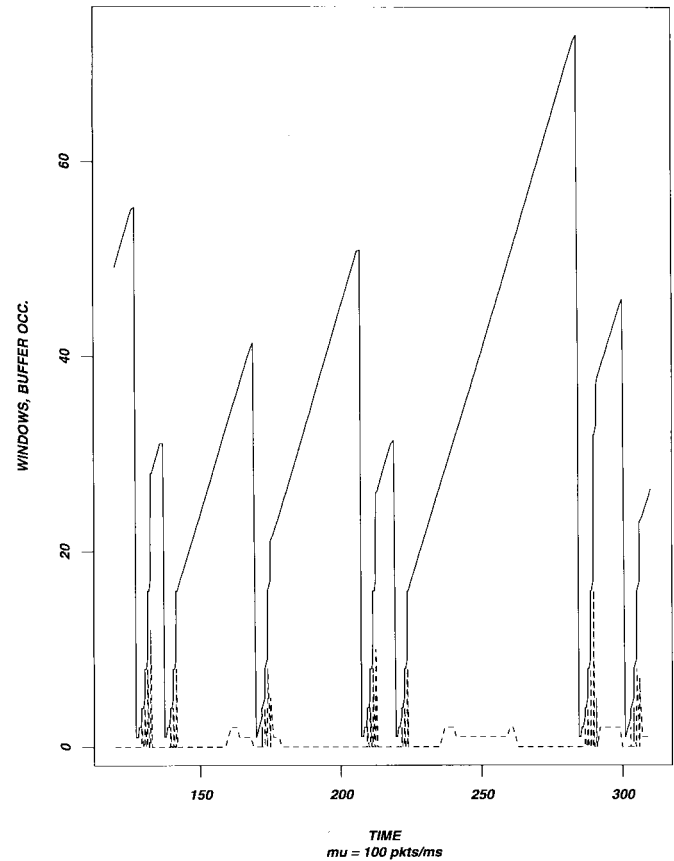


Fig. 3. Window and buffer evolution for single connection with high random losses. Prop. delay = 1 ms; $b = 0.8$; $q = 0.001$.

to further simplifications. Let $n(t, w)$ be the number of successful packets at time t , given that the cycle starts at time zero with slow start threshold w . Given w , the function $n(t, w)$ is the inverse of the function of $T(n, w)$, i.e., $t = T[n(t, w), w]$. However, this function is easier to compute, and indeed, has been computed in the previous section, since $dn/dt = \lambda(t) = \min\{W(t)/T, \mu\}$. Starting with (21), we change the summation to an integral and then change variables from n to t to obtain

$$\begin{aligned} \bar{\lambda} &\approx \int_n P[N = n] \frac{n}{T(n, w^*)} dn \\ &= \int_t P[N = n(t, w^*)] \frac{n(t, w)}{t} \lambda(t) dt. \end{aligned} \quad (27)$$

The analysis of the previous section can now be directly applied to evaluate this expression. The details are relegated to the Appendix. For either TCP-tahoe or TCP-reno, the preceding integral is simple, and we use the general-purpose tool Mathematica to compute it. The result is compared with simulations in Table III. We consider two values of normalized buffer size for a fixed value of bandwidth-delay product. Although the match in numerical results obtained from approximate analysis and simulation is not as good as in the previous section, the *qualitative* observations are identical. When the loss probability q is relatively high [an order of magnitude or more larger than $(\mu T)^{-2}$], the window size at the beginning of the congestion avoidance phase, and hence

TABLE III
LINK UTILIZATION AS A FUNCTION OF LOSS PROBABILITY q FOR
 $\mu = 100$, $\tau = 1$, AND TWO VALUES OF β ($\beta = 0.8$ AND $\beta = 0.2$)

q	Link Utilization (Analysis/Simulation)			
	$\beta = 0.8$		$\beta = 0.2$	
	Tahoe	Reno	Tahoe	Reno
10^{-1}	.021/.025	.035/.019	.021/.025	.035/.019
10^{-2}	.079/.098	.116/.108	.079/.095	.116/.108
10^{-3}	.291/.343	.369/.381	.291/.340	.369/.379
10^{-4}	.829/.861	.952/.911	.548/.627	.820/.795
10^{-5}	.940/.947	.994/.989	.647/.661	.865/.863
10^{-6}	.952/.953	.996/.994	.659/.656	.870/.870
0	.953/.954	.996/.994	.660/.664	.871/.870

TABLE IV
LINK UTILIZATION AS A FUNCTION OF $q(\mu\tau)^2$ FOR THREE
DIFFERENT VALUES OF BANDWIDTH-DELAY PRODUCT ($\beta = 0.8$)

$q(\mu\tau)^2$	Link Utilization (Tahoe/Reno)		
	$\mu = 100,$ $\tau = 1$	$\mu = 200,$ $\tau = 1$	$\mu = 200,$ $\tau = 2$
10	.343/.381	.363/.398	.377/.406
1	.861/.911	.887/.914	.906/.995
0	.954/.994	.969/.995	.981/.998

the throughput, is dominated by the effect of random losses. This throughput is much smaller than the lossless throughput and is insensitive to the value of β . As the loss probability decreases to $(\mu T)^{-2}$, the throughput gets closer to its lossless value, and the effect of the buffer size on throughput becomes apparent.

The strong dependence of link utilization on $q(\mu T)^2$ is made explicit in Table IV, where we show, for selected values of $q(\mu\tau)^2$ (note that $T = \tau + 1/\mu \approx \tau$ for large bandwidth-delay products), the link utilization for several different bandwidth-delay products. The normalized buffer size β is fixed at 0.8. The dependence on the bandwidth-delay product itself [for fixed $q(\mu\tau)^2$ and fixed β] is weaker, but there is a slight improvement in link utilization as $\mu\tau$ increases. For brevity, we show only the simulation results, but the analytical results exhibit the same qualitative features.

V. MULTIPLE TCP CONNECTIONS

In this section, we focus on the bias of the TCP window adjustment mechanism against connections with larger round-trip delays, and show that it is a fundamental property of the TCP window dynamics. Our simulations also reveal the disturbing impact of phase effects on TCP-reno. When simulating TCP-reno, therefore, we will consider a more detailed model of TCP connections going through an ATM switch with randomization in the cell discards when the buffer is full. These simulations serve a two-fold purpose: they enable us to focus on the window dynamics by eliminating the phase effects, and they verify that the analytical results, which are derived for a much simpler model, predict performance well for much more complicated system models.

An intuitive explanation of the bias against connections with longer round-trip times is provided in Section V-A. Section V-B contains a description of the ATM simulation model used for TCP-reno. Section V-C contains an approximate analysis of TCP-tahoe and TCP-reno for the simple system model we

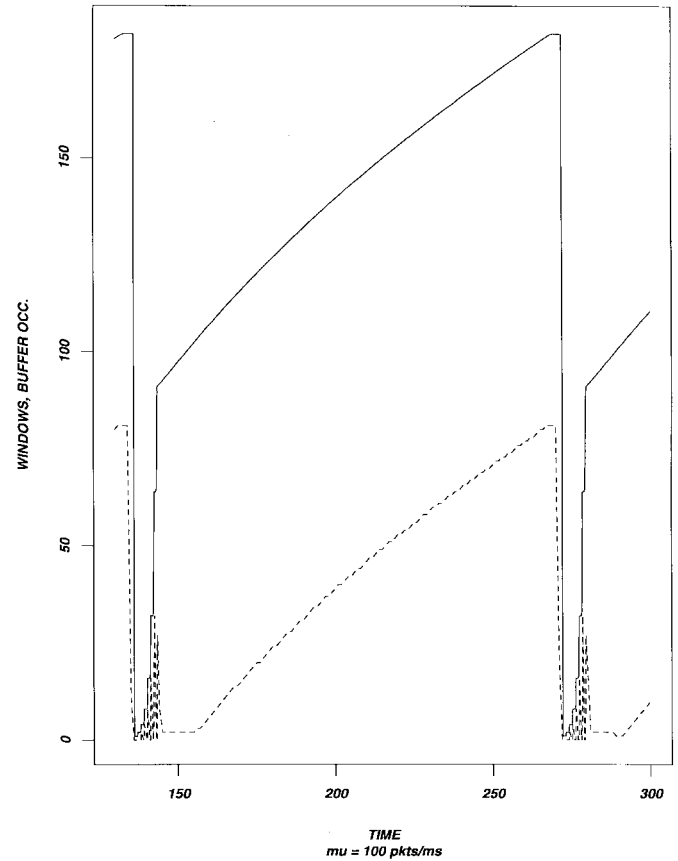


Fig. 4. Window evolution for connections with propagation delays 1 and 3. Prop. delay = 1 ms; $b = 0.8$.

have considered so far. Numerical results comparing the results of this analysis with simulations are given in Section V-D.

Throughout this section, we consider N TCP connections. Let $T_i = \tau_i + 1/\mu$ denote the round-trip delay for a packet of the i th connection if it arrives to an empty queue, where τ_i is the propagation delay for the connection. Fig. 4 shows the evolution of window sizes for two connections with propagation delays $\tau_1 = 1$ (solid line) and $\tau_2 = 3$ (dashed line). Note that the connection with larger propagation delay operates with a much smaller window size and consequently has a much lower relative throughput. The relative throughputs for the two connections are 0.817 and 0.114. For $\tau_1 = 1$ and $\tau_2 = 2$, the corresponding relative throughputs are 0.628 and 0.301. We observe empirically that the throughput is inversely proportional to τ_i^a , where $1 \leq a \leq 2$.

Another observation from our simulations is that the window evolution for the two connections become synchronized even when the two connections start at different times and despite the fact that the connections have different propagation delays. Synchronized window evolution has been previously reported in [24] for connections with equal propagation delays. The implication of this synchronized evolution is (as pointed out in [24]) that after a period of congestion the sum of connection window sizes may be too small to fill up the available network bandwidth, which results in throughput loss. For tractability, our approximate analysis is based on the assumption of synchronization. It appears difficult to deduce

analytically whether such synchronization would always occur, but the basic observation on which our analysis is based, which is that the window size grows more slowly for connections with higher round trip delays, should apply even if there were instances of evolution without synchronization.

A. Intuitive Explanation for the Bias

An approximate expression for the instantaneous throughput for connection i can be obtained using Little's law:

$$\lambda_i = \frac{W_i}{T_i + D_i} \quad (28)$$

where D_i is the typical waiting time in the queue for a packet from connection i . Note that although D_i is defined as a time-varying quantity, calling it "typical" implicitly assumes that the time variations are slow, as does the application of Little's law, which holds for averaged quantities. Clearly, if all connections had the same window size, and if D_i were small compared to T_i , then the throughput would be inversely proportional to T_i . Moreover, from (19), (20) in Section III-B, the growth in window size in the congestion avoidance phase is given by

$$\frac{dW_i}{dt} = \frac{\lambda_i}{W_i} = \frac{1}{T_i + D_i}, \quad 1 \leq i \leq N. \quad (29)$$

Assume now that the connections evolve in synchrony, i.e., that for all connections, the window sizes reach their maximum, and then drop down on detecting packet loss, at roughly the same time. Equation (29) shows that both the throughput and the growth in window size are inversely proportional to T_i (if the waiting times D_i are small). The maximum window size is therefore smaller for connections with larger propagation delays, so that in the next cycle, the initial window size during congestion avoidance is smaller for such connections. Thus, not only are the window sizes smaller throughout the cycle for such connections, but (28) implies that the throughput for such connections would be reduced by a further factor of roughly $1/T_i$. Thus, if there were no queueing delays, the average throughput for connection i could be expected to be inversely proportional to $1/T_i^2$. This unfairness is alleviated somewhat due to all connections suffering roughly equal queueing delays, which accounts for the empirical observation that the throughput is actually inversely proportional to $1/T_i^a$, where $1 \leq a < 2$. Thus, the bias against connections with larger propagation delays is a fundamental consequence of TCP dynamics. As we will see from the numerical results, this is alleviated somewhat when the buffer size is much larger than the bandwidth-delay product, since in that case the round-trip time is dominated by the queueing delay, which at high utilizations is roughly the same for all connections.

B. Simulation Model for TCP-reno

TCP-reno differs from TCP-tahoe mainly in its attempt to eliminate the relatively short slow start from the cyclical evolution by means of fast retransmit. Since its window dynamics are otherwise the same as that of TCP-tahoe, we would expect to see the same bias against connections with larger propagation delays in both versions of TCP. However,

our simulations for the simple model described in Section II reveal the following disturbing feature about TCP-reno: since the number of times the window is halved at the onset of congestion equals the number of lost packets, and since phase effects [9] can cause one connection to systematically lose a larger number of packets, it is possible that a connection gets almost completely shut out. We have observed this effect for two connections with $\tau_1 = 1$ and $\tau_2 = 3$, in which the first connection gets very small throughput even though it has smaller round-trip delay. This effect does not occur in TCP-tahoe because initiation of the slow start phase and the choice of the slow start threshold depends only on at least one packet getting lost, and not on the *number* of packets lost, at the onset of congestion.

In order to eliminate the possibility that phase effects are strictly an artifact of our simple model, we use a detailed simulation of multiple TCP-reno connections over an ATM system. Here, TCP sources are connected to routers with ATM interfaces. We use 576-byte packets which convert to 12 ATM cells at the router output. ACK packets convert to 2 ATM cells. The routers connect to the input ports of an ATM switch over links operating at DS-3 rates. The routers on the sending side assign the right virtual circuit (VCs) for transport through the ATM switch. The VCs in simulated system are configured such that the only bottlenecks are at the ATM switch output ports. The routers at the sending and receiving sides only act as nonbottlenecking hops.

The DS-3 links (using the Physical Layer Convergence Procedure over DS-3) typically operate in a slotted manner with slot duration equal to 10.4 μ s. However, on the input side to the switch, we do not use slotting on the DS-3 links, so that cell transmissions from the router do not have to wait for a slot boundary. The links connected to the switch output port operate in a slotted manner. The architecture and internal timing of a popular ATM switch are accurately simulated. The switch input ports (where the DS-3 links from the router terminate) are sampled every 2.8 μ s. A more detailed discussion of internal architecture is not relevant for the performance aspects studied in this paper. It is sufficient to point out that there is no queueing or blocking internal to the switch, and that contention occurs only at the switch output ports. Congestion is caused by routing many connections to the same output port. The output links from the switch terminate at routers which convert incoming cells to packets and pass the packet to appropriate receivers. Packets with lost cells are dropped at this point.

Mixing slotted and unslotted links is found not to eliminate phase effects. We therefore resort to randomization in buffer discards to eliminate phase effects. With randomization, when a cell arrives to a full buffer, instead of dropping the incoming cell, we choose the cell to be dropped randomly from a pool of cells, including the incoming cell and n cells in the tail of the buffer (n is chosen equal to the number of active TCP connections in our simulations).

C. Analysis and Numerical Results

We now attempt to develop analytical approximations for the throughput by refining the qualitative explanations given

earlier. The analysis is approximate for several reasons: i) an exact analysis of multiple connections with different propagation delays is not available even for fixed windows, which makes it necessary to approximate the queueing delays D_i in (28) and (29); ii) it is necessary to further approximate the result of step i) in order to solve analytically the differential equations arising from a continuous-time approximation to the window evolution; iii) our assumption of synchronized evolution for all connections does not hold exactly; and iv) our analysis is based on a “fixed point” argument, which assumes periodic behavior that may not hold exactly in the actual system.

In the case of TCP-reno, there is an additional approximation. The number of packets lost by a connection in a congestion episode determines what factor the window size is reduced by. This number varies over different congestion episodes for each connection, and it is difficult to take this variation into account in the analysis. Simulation results show that each connection typically loses about two packets on the average in each congestion episode. Losses occur when the buffer is full and one of the connections increases its window by one. When cells from this new packet arrive at the buffer, they typically cause losses to cells belonging to two packets (the tail of the packet arriving from the other connection and the head of the next). So on the average one would expect 3 packets to be lost per collision. In our example of two connections with 40- and 80-ms round trip times, once the buffer is full and a collision happens, the congestion episode lasts for 40 ms. During this period the 80 ms connection increases its window roughly 50% of the time, so that the average number of lost packets due to this increase is 1.5. Therefore 4.5 packets in total, or 2.25 packets per connection, are lost on the average per congestion episode. We will assume in our analysis that each connection loses *exactly* two packets in each congestion episode, so that the new window size is one-fourth of the window size at the onset of congestion.²

Note that the connections evolve in approximate synchrony. This is because whenever a collision happens due to a connection’s window increasing when the buffer is full, all cells arriving during the period over which the collision-causing packet arrives are dropped. With the assumption of infinite sources and cell jitters not more than a packet transmission time on an input link, all connections will transmit cells during the time over which the collision-causing packet arrives. Hence, all connections lose packets and halve their windows within the maximum round-trip time.

Assume now that the N TCP-tahoe connections evolve periodically and in synchrony as follows. In the congestion avoidance phase of each cycle, connection i goes from an initial window value of w_i to a value of $2^c w_i$, at which point packets are lost by each connection. Here, $c = 1$ for TCP-tahoe (assuming that each connection loses at least one packet), and $c = 2$ for TCP-reno (assuming that each connection loses exactly two packets). We therefore obtain that the window

²The preceding assumption is not necessary for TCP-tahoe. As long as each connection loses at least one packet in a congestion episode, all the window sizes drop to one, and the slow start threshold for each connection is set at half of its current window size.

size at the beginning of the congestion avoidance phase is again w_i (ignoring the slow start phase for the next cycle in the case of TCP-tahoe). We would like to find the value of the w_i such that the preceding “fixed point” behavior holds, and use the resulting window evolution to estimate the throughput obtained by each connection.

In order to solve (29), we must estimate the queueing delays D_i . At any instant of time, the $D_i(t)$ can be taken to be the average delays for an analogous system with *fixed* windows $W_i(t)$. If the link utilization is known to be 100%, and if it is assumed that $D_i \equiv D$ for all i (this has been validated by simulations for fixed windows), then, from (28), the delay D must satisfy the following equation:

$$\sum_{i=1}^N \frac{W_i}{T_i + D} = \mu, \quad \text{for 100\% utilization.} \quad (30)$$

At present, we do not know how to compute the delay when the utilization is *less than* 100%. In fact, we do not even have a criterion for when the utilization is 100% (we have derived necessary and sufficient conditions, but these do not coincide). Finally, solving (29) analytically when the D_i are given by the implicit expression (30) appears difficult. We therefore use the following simple approximations. We divide the evolution in the cycle into two phases:

$$\begin{aligned} \text{Phase A:} & \quad \sum_{i=1}^N \frac{W_i}{T_i} < \mu \\ \text{Phase B:} & \quad \sum_{i=1}^N \frac{W_i}{T_i} \geq \mu. \end{aligned}$$

The condition for Phase B can be shown to be a necessary (but not sufficient) condition for 100% utilization. If the buffer size is large compared to the bandwidth-delay product, Phase A may not occur at all, since the link may be fully utilized even when the window sizes are reduced following a congestion episode. We must therefore consider two different cases in the analysis.

Case 1—Small Buffers: Assume that the link is not always fully utilized during congestion avoidance, and that the queueing delays for all connections are the same, $D_i \equiv D$, where D is given by

$$D \approx \begin{cases} 0, & \text{in Phase A} \\ \frac{B}{2\mu}, & \text{in Phase B.} \end{cases} \quad (31)$$

The motivation for the preceding approximation is as follows. Since the utilization is known to be less than 100% in Phase A, we assume that there is no waiting in queue in this phase. Phase B, on the other hand, starts from a delay of zero, and ends with a delay of B/μ , assuming that the typical queue length is approximately the same as the maximum queue length toward the end of Phase B (recall that Phase B ends when the latter exceeds B). We therefore use $B/(2\mu)$, the “average” delay value during Phase B, as the value of the queueing delay used in (29) for this phase.

Substituting (31) in (29), we obtain

$$\frac{dW_i}{dt} = \begin{cases} \frac{1}{T_i}, & \text{in Phase A} \\ \frac{1}{T_i + \frac{B}{2\mu}}, & \text{in Phase B.} \end{cases} \quad (32)$$

Starting from an initial value of w_i , we obtain that

$$W_i(t) = w_i + \frac{t}{T_i}, \quad 0 \leq t \leq t_A \quad (33)$$

$$W_i(t + t_A) = W_i(t_A) + \frac{t}{T_i + \frac{B}{2\mu}}, \quad 0 \leq t \leq t_B \quad (34)$$

where t_A, t_B denote the durations of Phase A and B, respectively. These quantities are obtained as (linear) functions of the w_i as follows. We have

$$\sum_{i=1}^N \frac{W_i(t_A)}{T_i} = \mu,$$

which yields, using (33), that

$$t_A = \frac{\mu - \sum_{i=1}^N \frac{w_i}{T_i}}{\sum_{i=1}^N \frac{1}{T_i^2}}. \quad (35)$$

We now compute the duration of the second phase. This phase lasts until the delay grows to B/μ , at which point the current cycle terminates. Since the link is fully utilized at this time, using (28), we obtain that the duration t_B of this phase is given by the equation

$$\sum_{i=1}^N \frac{W_i(t_A + t_B)}{T_i + \frac{B}{\mu}} = \mu. \quad (36)$$

Substituting from (33) to (35), we obtain t_B as a linear function of the w_i . Since the window size at the end of the second phase is 2^c times the starting window size for the next cycle, invoking the fixed point condition gives that

$$W_i(t_A + t_B) = 2^c w_i, \quad 1 \leq i \leq N, \quad (37)$$

These are N linear equations in the N unknowns, w_1, \dots, w_N . While we ignored the slow start phase in computing the slow start thresholds w_i , we can now estimate the duration and number of packets for this period as $t_{ss,i} = T_i \log_2 w_i$ and $n_{ss,i} = w_i$, respectively. This yields the following estimate for the average throughput for the i th connection:

$$\bar{\lambda}_i = \frac{\int_0^{t_A} \frac{W_i(t)}{T_i} dt + \int_0^{t_B} \frac{W_i(t + t_A)}{T_i + \frac{B}{2\mu}} dt}{t_A + t_B} \quad \text{reno}$$

$$\bar{\lambda}_i = \frac{w_i + \int_0^{t_A} \frac{W_i(t)}{T_i} dt + \int_0^{t_B} \frac{W_i(t + t_A)}{T_i + \frac{B}{2\mu}} dt}{T_i \log_2 w_i + t_A + t_B}. \quad \text{tahoe}$$

(While the preceding expressions are similar, note that the value of c used in (37) to compute the times t_A and t_B is $c = 1$ for TCP-tahoe and $c = 2$ for TCP-reno.)

Case 2—Large Buffers: The analysis of Case 1 does not apply when the buffer size is large enough that the link is fully utilized throughout congestion avoidance, despite the window size reduction following the onset of congestion. In fact, the analysis in Case 1 yields that t_A is negative, which therefore provides a simple criterion as to when to apply the analysis in Case 2 which is described in the following.

Assuming that the link is always fully utilized during congestion avoidance, let D_{\min} be the (unknown) minimum value of the queuing delay seen at the beginning of congestion avoidance. Using the average of this and the maximum queuing delay B/μ to linearly approximate the window evolution described by (29), we obtain that

$$W_i(t) = w_i + \frac{t}{T_i + \bar{D}}, \quad 0 \leq t \leq t_B \quad (38)$$

where w_i is the window size of the i th connection at the beginning of congestion avoidance. Using the fixed point argument,

$$W_i(t_B) = 2^c w_i. \quad (39)$$

It follows from (45) and (46) that

$$\frac{w_i}{w_1} = \frac{T_1 + \bar{D}}{T_i + \bar{D}}. \quad (40)$$

Since the link is fully utilized, the following equations must hold:

$$\sum_{i=1}^N \frac{w_i}{T_i + D_{\min}} = \mu \quad (41)$$

$$\sum_{i=1}^N \frac{2^c w_i}{T_i + \frac{B}{\mu}} = \mu. \quad (42)$$

Using (40)–(42), we can eliminate the w_i and solve for D_{\min} . We can substitute back to solve for the w_i and the duration t_B of the congestion avoidance phase. The throughputs are given by

$$\bar{\lambda}_i = \frac{\int_0^{t_B} \frac{W_i(t)}{T_i + \bar{D}} dt}{t_B} \quad \text{TCP-reno}$$

$$\bar{\lambda}_i = \frac{w_i + \int_0^{t_B} \frac{W_i(t)}{T_i + \bar{D}} dt}{T_i \log_2 w_i + t_B} \quad \text{TCP-tahoe}$$

TABLE V
RELATIVE THROUGHPUTS FOR TWO TCP-TAHOE CONNECTIONS WITH
DIFFERENT PROPAGATION DELAYS FOR $\mu = 100$ AND BUFFER SIZE $B = 80$

Prop. Delays		Relative Throughputs			
		α_1		α_2	
τ_1	τ_2	Analysis	Simulation	Analysis	Simulation
2	1	.232	.301	.704	.628
3	1	.133	.114	.809	.817

TABLE VI
RELATIVE THROUGHPUTS FOR TWO CONNECTIONS WITH $\tau_1 = 80$
MS, $\tau_2 = 40$ MS, $\mu = 96\,000$ CELLS/S (DS-3 LINK RATE)

Buffer size (cells)	Relative Throughputs			
	α_1		α_2	
	Analysis	Simulation	Analysis	Simulation
256	.1324	.0892	.5285	.5268
2304	.1862	.1385	.6639	.7033
5760	.2539	.2316	.7042	.6966
8064	.2876	.2083	.6918	.7410
14400	.3520	.3388	.6389	.6217
23040	.3941	.3199	.6007	.6318

D. Numerical Results

For TCP-tahoe, analysis and simulation results are compared for the simple model described in Section II. Table V gives the relative throughputs $\alpha_i = \bar{\lambda}_i/\mu$ for two connections with different propagation delays. The agreement between analysis and simulation is seen to be quite good.

For TCP-reno, simulation results obtained for the ATM system described in Section V-C are compared with analytical results for the simple model of Section II. Table VI gives the relative throughputs $\alpha_i = \bar{\lambda}_i/\mu$ for two connections with different propagation delays. As described earlier, the analytical results are obtained assuming that exactly 2 losses occur per congestion episode, so that $c = 2$ in (37) and (39). Despite this coarse characterization of the loss behavior at the onset of congestion, and despite all the other simplifying approximations we have made, the agreement between analysis and simulation is quite good for all values of the buffer size.

VI. CONCLUSIONS

The results stated in the introduction and developed in the rest of the paper imply the following observations (some of which are speculations regarding suitable directions for future research) regarding TCP in high-bandwidth delay product networks.

- 1) The use of cumulative acknowledgment in TCP motivates the TCP-tahoe feature of reducing the window size to one after a loss in order to avoid a burst of packets when the retransmission gets through. This in turn makes an exponential increase of window size in slow start necessary, especially for high bandwidth-delay networks, in order obtain nontrivial link utilizations. On the other hand, this exponential increase causes bursty traffic, which, if the buffer size is smaller than 1/3 of the bandwidth-delay product, causes buffer overflow and a second slow start phase, leading to a lower throughput. TCP-reno tries to avoid this phenomenon by cutting the window to half when it detects a loss. While this does provide better throughput under ideal conditions,

TCP-reno in its present form is too vulnerable to phase effects and multiple packet losses to be a replacement for TCP-tahoe. The basic problem with TCP-reno is that there can be multiple window cutbacks due to a single congestion episode, and that multiple losses can lead to a timeout (which in practice can lead to significant throughput loss if coarse timeouts are used). A recently proposed version of TCP (TCP-Vegas) [4] attempts to address this problem, in addition to other proposed changes such as more sophisticated processing of round-trip time estimates. A detailed discussion of TCP-Vegas is beyond the scope of this paper. However, in our opinion, in order to significantly improve upon the present versions of TCP, it is necessary to avoid the drastic window reductions in both TCP-tahoe and TCP-reno except when there is sustained congestion (which would cause multiple losses). One possible means of handling isolated losses without changing the window size is by using some form of selective acknowledgment (in which case TCP dynamics should be reworked to take advantage of such an added feature). Pending development of a satisfactory replacement, we recommend the use of TCP-tahoe (together with network level controls to optimize its performance), since it is far more robust than TCP-reno.

- 2) TCPs vulnerability to random loss makes it difficult to multiplex data traffic with real-time traffic with a rapidly time-varying rate, especially if both kinds of traffic share the same buffer, as is the case with most current networks. In emerging ATM networks, however, TCP connections might be supported over the UBR or ABR traffic classes, which would typically be buffered separately from higher priority variable bit rate (VBR) or constant bit rate (CBR) traffic. If the latter traffic classes have higher priority, TCP connections would see a time-varying link capacity left over after VBR and CBR traffic has been served. These variations could lead to "random loss," which we have shown can seriously impact performance. However, if there is sufficient buffering to absorb these variations and to keep the end-to-end packet loss probability below the inverse square of the bandwidth-delay product, then our results also imply that TCP performance would not be seriously affected. The results of this paper could therefore form a basis for further investigation as to the buffer sizes required to support TCP over ABR or UBR service classes in ATM. Indeed, we have applied the results here to the design of wireless-wireline interfaces to "hide" the time variations in a wireless channel from TCP/IP [5], so that the random loss seen by TCP is small enough to maintain a high throughput.
- 3) The unfairness of TCP toward connections with higher propagation delays could cause performance problems when multiplexing short and long-haul traffic on WANs. For guaranteed performance in highly utilized networks, each TCP connection should be given reserved buffer and bandwidth resources throughout the network. Typically, the resource allocation would be determined at

connection set-up and enforced at switches and routers using per connection queueing [6], [21]. Since administering the resources allocated to *every* best effort connection may be excessively expensive, a more feasible alternative might be to allocate and administer resources for an entire traffic class. In such a situation, the unfairness we have pointed out would persist if TCP were supported over the ATM UBR traffic class. However, if TCP is supported over the ABR traffic class, the time-varying rate available to each connection is determined at the network level and is administered at the source, so that different TCP connections should be isolated from each other to a large extent even if they share the same network buffers.

- 4) In addition to causing vulnerability to random loss, the fact that loss is the sole means of feedback used by TCP leads to excessive delays. This is because, in networks with high utilizations, the window size for a TCP connection would keep increasing after the bottleneck link is fully utilized, until in fact there is a buffer overflow leading to a loss. The delay and loss performance would improve significantly if we instead used a scheme that tries to maintain a window size which is just large enough to achieve a high link utilization. A scheme such as DECBIT [22] attempts to do this using explicit feedback from the switches, and similar schemes are worth pursuing, especially because Explicit Congestion Notification is incorporated as an option for ATM networks. Note, however, that the DECBIT scheme in particular shares with TCP the problem of unfairness toward connections with longer propagation delays. Another possibility is a more sophisticated processing of round-trip time estimates similar to the approach taken in [18], [19]. This is certainly attractive, since it avoids the need for explicit feedback. However, if the round-trip delays can change substantially without changes in the load on the path of the connection (e.g., because processing delays at nodes depend on the load on the operating system, or because of delays due to handoffs for mobile computing applications), then adaptation based on delay processing might be less robust than adaptation based on loss or explicit feedback. In addition, if different connections are not isolated from each other in terms of their use of bandwidth and buffering in the network, then a connection that is more aggressive about obtaining bandwidth by increasing its rate until there is a loss would be at an advantage over connections that process round-trip delays to avoid congestion. Thus, changes at the transport layer must either be adopted universally, or must go hand in hand with network layer controls that guard against greedy connections.

In summary, while we have identified several shortcomings of TCP, we have also mentioned possible means of obtaining good performance via network level solutions, such as isolating connections from each other and providing enough buffering to hide excessively fast time variations in available link capacity from TCP. An important topic for future research

is, for each context of interest, to translate these into specific recommendations, and to provide systematic design techniques for arriving at these recommendations. Especially interesting is the question of how best to support TCP over the ATM ABR and UBR service classes, since that involves adaptation at both the network and transport layers. Another important area for future research is the development of an alternative dynamic window mechanism which addresses some of the shortcomings of TCP while preserving its decentralized nature. Possible improvements might be better congestion avoidance via more sophisticated processing of round-trip delay estimates, and the use of selective acknowledgments to improve the performance in the presence of random loss.

APPENDIX

We give the details of the approximation for the throughput with loss, taking as our starting point (27) in Section IV.

We consider only the case $\beta \leq 1$. Considering TCP-reno first, with the threshold at the beginning of the cycle given by $W_0 = w^*$ given by (26) instead of $W_{\max}/2$. Referring to the analysis in Section III of the congestion avoidance phase, the instantaneous throughput $\lambda(t)$ is given by

$$\lambda(t) = \begin{cases} \frac{W_0 + \frac{t}{T}}{T}, & 0 \leq t \leq t_A \\ \mu, & t_A \leq t \leq t_A + t_B \end{cases} \quad (43)$$

and the number of successful packets by time t is there is no random loss is given by

$$n(t) = \int_0^t \lambda(t) dt. \quad (44)$$

Given the simple form of (43), it is easy to compute $n(t)$ explicitly using (44). Substituting into (27), and using (18), we obtain

$$\bar{\lambda} \approx \int_0^{t_A+t_B} q(1-q)^{\frac{n(t)}{t}} \lambda(t) dt + (1-q)^{N_{\max}} \frac{N_{\max}}{t_A+t_B}, \quad \text{TCP-reno.} \quad (45)$$

This is the expression that is plugged into Mathematica to generate the desired results.

For TCP-tahoe, according to our approximation, the threshold for the slow start phase is now given by $W_t = w^*$. The case-wise analysis of Section III-A for this value of W_t yields t_{ss} , n_{ss} , and W_0 . We assume for simplicity that there is no random loss in the slow start phase(s), since relatively few packets are transmitted in this phase. As for TCP-reno, the instantaneous throughput in the congestion avoidance phase is given by (43) (it is convenient to take $t = 0$ to correspond to the beginning of congestion avoidance rather than to the beginning of the cycle), and the number of successful packets by time t is given by $n(t) = n_{ss} + \int_0^t \lambda(t) dt$. A simple modification of (52) which includes the duration t_{ss} of the slow start phase yields the following approximation for the

average throughput:

$$\bar{\lambda} \approx \int_0^{t_A+t_B} q(1-q)^{n(t)} \frac{n(t)}{t+t_{ss}} \lambda(t) dt + (1-q)^{N_{\max}} \frac{N_{\max}}{t_{ss}+t_A+t_B}, \quad \text{TCP-tahoe.} \quad (46)$$

ACKNOWLEDGMENT

The authors would like to thank T. Ott for getting them started on this problem.

REFERENCES

- [1] *ATM Forum Traffic Management Specification Version 4.0, Draft Specification ATM Forum/95-0013R11*, ATM Forum, Mar. 1996.
- [2] J. Bolot and A. U. Shankar, "Dynamical behavior of rate-based flow control mechanisms," *Comp. Comm. Rev.*, pp. 35–49, Apr. 1990.
- [3] J. Bolot, "End-to-end packet delay and loss behavior in the Internet," in *Proc. ACM SIGCOMM'93*.
- [4] L. S. Brakmo, S. W. O'Malley, and L. L. Peterson, "TCP Vegas: New techniques for congestion detection and avoidance," in *Proc. ACM Sigcomm*, Aug. 1994.
- [5] H. Chaskar, T. V. Lakshman, and U. Madhow, "On the design of interfaces for TCP/IP over wireless," in *Proc. IEEE Milcom'96*.
- [6] A. Demers, S. Keshav, and S. Shenker, "Analysis and simulation of a fair queueing algorithm," in *Proc. ACM SIGCOMM'89*.
- [7] K. W. Fendick, D. Mitra, I. Mitrani, M. A. Rodrigues, J. B. Seery, and A. Weiss, "An approach to high performance, high speed data networks," *IEEE Comm. Mag.*, pp. 74–82, Oct. 1991.
- [8] S. Floyd, "Connections with multiple congested gateways in packet-switched networks, Part 1: One-way traffic," *Comp. Comm. Rev.*, vol. 21, no. 5, pp. 30–47, Oct. 1991.
- [9] S. Floyd and V. Jacobson, "On traffic phase effects in packet-switched gateways," *Internetworking: Research and Experience*, vol. 3, no. 3, pp. 115–156, Sept. 1992. (An earlier version of this paper appeared in *Comp. Comm. Review*, vol. 21, no. 2, Apr. 1991.)
- [10] —, "Random early detection gateways for congestion avoidance," *IEEE/ACM Trans. Networking*, vol. 1, no. 4, pp. 397–413, Aug. 1993.
- [11] V. Jacobson, "Congestion avoidance and control," in *Proc. ACM SIGCOMM'88*, pp. 314–329.
- [12] —, "Modified TCP congestion avoidance algorithm," message to end2end-interest mailing list, Apr. 1990, URL <ftp://ftp.ee.lbl.gov/email/vanj.90apr30.txt>.
- [13] —, "Berkeley TCP evolution from 4.3-tahoe to 4.3-reno," in *Proc. 18th Internet Engineering Task Force*, Vancouver, Aug. 1990.
- [14] V. Jacobson, R. Braden, and D. Borman, "TCP extensions for high performance," RFC (request for comment) 1323, May 1992.
- [15] P. Karn and C. Partridge, "Improving round-trip time estimates in reliable transport protocols," *ACM Trans. Comp. Sys.*, vol. 9, no. 4, pp. 364–373, Nov. 1991.
- [16] T. V. Lakshman, U. Madhow, and B. Suter, "Window-based error recovery and flow control with a slow acknowledgment channel: A study of TCP/IP performance," in *Proc. IEEE Infocom 1997*.
- [17] B. Makrucki, "On the performance of submitting excess traffic to ATM networks," in *Proc. Globecom 1991*, Dec. 1991, pp. 281–288.
- [18] D. Mitra, "Asymptotically optimal design of congestion control for high speed data networks," *IEEE Trans. Commun.*, vol. 40, no. 2, pp. 301–311, Feb. 1992.
- [19] D. Mitra and J. B. Seery, "Dynamic adaptive windows for high speed data networks with multiple paths and propagation delays," *Computer Networks and ISDN Systems*, vol. 25, pp. 663–679, 1993.
- [20] A. Mukherjee and J. C. Strikwerda, "Analysis of dynamic congestion control protocols—A Fokker–Planck approximation," in *Proc. ACM SIGCOMM'91*, pp. 159–169.

- [21] A. K. Parekh and R. G. Gallager, "A generalized processor sharing approach to flow control in integrated services networks—The multiple node case," in *Proc. IEEE Infocom'93*.
- [22] K. K. Ramakrishnan and R. Jain, "A binary feedback scheme for congestion avoidance in computer networks with a connectionless network layer," in *Proc. ACM SIGCOMM'88*, pp. 303–313.
- [23] S. Shenker, "A theoretical analysis of feedback flow control," in *Proc. ACM SIGCOMM'90*, pp. 156–165.
- [24] S. Shenker, L. Zhang, and D. D. Clark, "Some observations on the dynamics of a congestion control algorithm," *Comp. Comm. Review*, pp. 30–39, Oct. 1990.
- [25] G. R. Wright and W. R. Stevens, *TCP/IP Illustrated, Volume 2, The Implementation*. Reading, MA: Addison Wesley, 1995.
- [26] L. Zhang, "A new architecture for packet switching network protocols," Ph.D. dissertation, MIT Lab. Comput. Sci., Cambridge, MA, 1989.
- [27] L. Zhang, S. Shenker, and D. D. Clark, "Observations on the dynamics of a congestion control algorithm: The effects of two-way traffic," in *Proc. ACM SIGCOMM'91*, pp. 133–147.



T. V. Lakshman (M'86) received the M.S. and Ph.D. degrees in computer science from the University of Maryland, College Park, in 1984 and 1986. Prior to that he received a Master's degree from the Department of Physics, Indian Institute of Science, Bangalore, India.

From 1986 to 1995, he was at Bellcore where he was most recently a Senior Research Scientist and Technical Project Manager in the Information Networking Research Laboratory. He is currently with the High Speed Networks Research Department at Bell Labs. His recent research has been in issues related to traffic characterization and provision of quality of service for video services, end-to-end flow control in high-speed networks, traffic shaping and policing, ATM switching, and parallel architectures for fast signaling and connection-management in high-speed networks. His current research interests are in the areas of high-speed networking, distributed computing, and multimedia systems.

Dr. Lakshman is a co-recipient of the 1995 ACM Sigmetrics/Performance Conference Outstanding Paper Award. He is an Editor of the IEEE/ACM TRANSACTIONS ON NETWORKING.



Upamanyu Madhow (SM'96) received the bachelor's degree in electrical engineering from the Indian Institute of Technology, Kanpur, India, in 1985. He received the M.S. and Ph.D. degrees in electrical engineering from the University of Illinois, Urbana-Champaign, in 1987 and 1990, respectively.

From August 1990 to July 1991, he was a Visiting Assistant Professor at the University of Illinois. From August 1991 to July 1994, he was a research scientist at Bell Communications Research, Morristown, NJ. Since August 1994, he has been an Assistant Professor with the Department of Electrical and Computer Engineering at the University of Illinois, Urbana-Champaign. His current research interests are in communication systems and networking, with current emphasis in wireless communications and high speed wide area networks.

Dr. Madhow is a recipient of the NSF CAREER award. He was awarded the President of India Gold Medal for graduating at the top of his undergraduate class.