

2000-07-03

TCP Congestion Control and Heavy Tails

Guo, Liang; Crovella, Mark; Matta, Ibrahim. "TCP Congestion Control and Heavy Tails", Technical Report BUCS-2000-017, Computer Science Department, Boston University, July 3, 2000. [Available from: <http://hdl.handle.net/2144/3774>]

<https://hdl.handle.net/2144/3774>

Boston University

TCP Congestion Control and Heavy Tails

LIANG GUO

MARK CROVELLA

IBRAHIM MATTA

Computer Science Department
 Boston University
 Boston, MA 02215

{guo1, crovella, matta}@cs.bu.edu

Abstract—Long-range dependence has been observed in many recent Internet traffic measurements. Previous studies have shown that there is a close relationship between heavy-tailed distribution of various traffic parameters and the long-range dependent property. In this paper, we use a simple Markov chain model to argue that when the loss rate is relatively high, TCP’s adaptive congestion control mechanism indeed generates traffic with heavy-tailed OFF, or idle, periods, and therefore introduces long-range dependence into the overall traffic. Moreover, the degree of such long-range dependence, measured by the Hurst parameter, increases as the loss rate increases, agreeing with many previous measurement-based studies. In addition, we observe that more variable initial retransmission timeout values for different packets introduces more variable packet inter-arrival times, which increases the burstiness of the overall traffic. Finally, we show that high loss conditions can lead to a heavy-tailed distribution of transmission times even for constant-sized files. This means that file size variability need not be the only cause of heavy-tailed variability in transmission durations.

Keywords— Congestion Control, Long-Range Dependence, Self-Similarity.

I. INTRODUCTION

SELF-SIMILARITY has been observed in a large number of Internet traffic measurements (for example, [1, 2, 3]). Self-similarity refers to the condition in which the second-order statistics of traffic decay very slowly with increasing the aggregation level, compared to traditional Markovian (memory-less) models. Self-similarity is closely associated with the phenomenon of heavy-tailed distributions, which are distributions whose tails decline via a power law with small exponent (less than 2). The presence of heavy tails in lengths of individual flows can be shown to induce self-similarity in network traffic [4]. Heavy-tailed properties have been found in file sizes and user thinking time [1], flow (session) duration [5], as well as packet inter-arrival time [4] distributions in the Internet.

To date, the best-accepted hypothesis for the genesis of self-similarity on time-scales from seconds to an hour is the heavy-tailed distribution that is typical of flow lengths

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

in the Internet [1, 4]—the majority of flows are found to be very short, many are long, and some are very long. However, a number of studies have shown that interesting scaling properties can arise even when flow lengths are *not* highly variable [6, 7]. Furthermore, although heavy-tailed flow lengths are commonly associated with heavy-tailed file sizes, the authors in [8] and [9] find *no* strong correlation between file sizes and transmission times (even though both show heavy tails).

In this paper, we study the relationship between TCP’s congestion control mechanism and traffic self-similarity under certain network conditions. By means of Markovian analysis and NS simulator [10], we confirm that TCP flows can show a heavy-tailed distribution in packet inter-arrivals under certain conditions. We show that larger loss rates bring longer range dependence into the traffic. This is consistent with previous measurement-based studies [3]. In addition, we also show in this paper that, because short connections do not have enough packet samples to accurately estimate the time they need to wait for acknowledgments, the conservatively chosen value can introduce very large gaps between packets, thus making the traffic generated by (short) TCP connections even more bursty.

Heavy-tailed Distribution, Self-Similarity and Long-Range Dependence: Here we briefly review some concepts related to fractal property. A more detailed description can be found in [11, 12] and references therein.

Most of the distributions shown in this paper have the property of being *heavy-tailed*. We say a distribution is heavy-tailed if the asymptotic shape of the distribution is power-law with exponent less than 2 regardless of the behavior of the distribution for small values of the random variable, i.e.,

$$P[X \geq x] \sim x^{-\alpha}, \quad \text{as } x \rightarrow \infty, 0 < \alpha < 2$$

The reason that such distributions are called heavy-tailed is that, compared to those more commonly used distributions such as exponential and normal distributions, a random variable that follows a heavy-tailed distribution can give rise to extremely large values with non-negligible probability. As a consequence, such random variable

shows infinite variance when $\alpha < 2$.

It has been shown in [4] that the aggregation of i.i.d. ON/OFF processes produces *self-similar* time-series, if either ON or OFF periods of each process follow a heavy-tailed distribution. We say a time-series is (asymptotically) self-similar if the autocorrelation function of the new time-series produced by aggregating the original time-series is (asymptotically) equal to the original autocorrelation function. That is, given a stationary time-series $X = (X_t : t = 0, 1, 2, \dots)$, we define the m -aggregated series $X^{(m)} = (X_k^{(m)} : k = 1, 2, 3, \dots)$ by summing the original series X over non-overlapping blocks of size m . Then if X is self-similar, it has the same autocorrelation function $r(k) = E[(X_t - \bar{X})(X_{t+k} - \bar{X})]$ as the series $X^{(m)}$ for all m , where $\bar{X} = E[X]$.

As a result, self-similar processes show *long-range dependence*. A process with long-range dependence has an autocorrelation function $r(k) \sim k^{-\beta}$ as $k \rightarrow \infty$, where $0 < \beta < 1$. Thus the autocorrelation function of such a process decays hyperbolically (as compared to the exponential decay exhibited by traditional Markovian traffic models).

One of the attractive features of using self-similar models for time-series, when appropriate, is that the degree of self-similarity of a series is expressed using only a single parameter called the *Hurst* parameter $H = 1 - \beta/2$. For self-similar series, $1/2 < H < 1$. As $H \rightarrow 1$, the degree of self-similarity increases. Thus the fundamental test for self-similarity of a series reduces to the question of whether H is significantly different from $1/2$.

Related Work: There has been a large body of work attempting to explain the causes of heavy-tails and self-similarity in Internet traffic. Broadly, they either attribute the causes to *application/user-level* variability (e.g., [13],[14],[1]) or to *system/network-level* complexity (e.g., [15],[7],[16],[17]).

Willinger *et al.* [13] examined Ethernet traffic at the packet level, identified flows between individual source/destination pairs, and showed that transmission and idle times for those flows are heavy-tailed. Paxson and Floyd [14] traced Internet traffic and observed that burst sizes in FTP transfers, and TELNET packet inter-arrival times (appearance of “keystroke”) show heavy-tailed distributions. In [1], Crovella and Bestavros examine Web traffic. They observe that file size distributions in Unix systems as well as in Web databases are heavy-tailed. In addition, heavy-tailed distributions are found in HTTP session time and request inter-arrival time (user “thinking” time). They propose that these heavy-tailed distributions might be the primary causes of the self-similarity in Web

traffic.

Other studies have argued that the chaotic nature of network protocols and variability of system conditions may also contribute to the self-similarity in traffic, especially at smaller time scales. In [15], Deane *et al.* analyze token-passing networks and attribute self-similarity in traffic to protocol interactions between sources and destinations. They show that even in a deterministic network, protocol dynamics can generate chaotic behavior, which leads to self-similarity in network traffic. A recent work by Veres *et al.* [7] illustrates by simulation that under severe network conditions, TCP congestion control protocol indeed shows chaotic nature, and starts to generate self-similar traffic. However, they only show evidence of such chaotic nature, and don’t explain why it shows up when the network is highly congested. In [16], Peha uses simulation to argue that packet retransmission and congestion control mechanisms could cause self-similarity when congestion does happen in the network. However, no theoretical evidence is given in that paper. More recently, Veres *et al.* [17] observe from real measurements that short TCP connections produce self-similar traffic, which they attribute to the reaction of TCP congestion control to the self-similar background traffic.

Our Contribution: We also conjecture that protocol-induced chaos could be a contributor to network traffic self-similarity. Our work is different from previous work in that we seek to discover the causes of protocol-induced chaos based on analytical arguments. More specifically, we illustrate in this paper that the exponential-backoff algorithm, used by TCP’s congestion control mechanism under severe network congestion conditions, can cause chaotic behavior. Our analytical model encompasses general loss conditions, which may be due to contention among several TCP connections (as in [7]), or due to self-similar cross-traffic (as in [17]). Thus, we view our work here as a first step toward analytically grounded evidence and reasons for self-similarity in TCP packet flows.

The paper is organized as follows. In Section II we briefly review TCP’s congestion control mechanism. In Section III we propose a Markovian model to describe its behavior under certain network conditions and analytically explain why the inter-arrival times of TCP packets follow a heavy-tailed distribution. We use simulation to confirm our analysis in Section IV and show that the traffic generated by TCP connections is indeed asymptotically self-similar. We discuss extensions to our analysis and future work in Section V.

II. CONSERVATIVE NATURE OF TCP AND ITS PATHOLOGICAL EFFECT

The objective of TCP congestion control is for each source to determine how much capacity is available in the network, so that it knows what rate it can safely send at. Therefore, the data transfer of TCP starts from a stage, called *slow-start*, in which TCP tries to increase its sending rate exponentially, until it encounters the first loss. At this point, TCP interprets packet loss as an indication of reaching the upper limit of the available bandwidth of the bottleneck link. Thus, from this point on (or following another slow-start period, depending on the implementation of TCP), it switches to another stage, called *congestion-avoidance*, in which TCP employs the *Additive Increase, Multiplicative Decrease (AIMD)* mechanism to slowly adapt to the available bandwidth.

Another important stage of TCP congestion control happens when the network is heavily loaded, during which some of the TCP connections should keep silent so as to clear congestion. This stage, referred to in this paper as *exponential-backoff*, is the regime that Karn's algorithm [18] deals with. In this regime, when TCP does not receive an acknowledgment for a packet after some timeout period, it assumes that this packet is lost, and then retransmits that packet and doubles its retransmission timeout value (RTO) for detecting packet loss. This process continues until the packet is successfully transmitted and acknowledged, up to some upper limit (usually 64 times the smallest timeout). Essentially, TCP tries to clear congestion by cutting its sending rate in half (or exponentially decreasing its rate).

Figure 1 shows a schematic view of the TCP congestion window behavior at different stages (black points on the top indicate packet losses). The value of the congestion window is proportional to the sending rate, which is roughly equal to the size of the window divided by the round-trip time (RTT). To simplify our analysis, we focus on the behavior of TCP Tahoe, and we assume that the receiver has an unlimited buffer, so there's no TCP flow control (the upper limit on the window size). We should point out that all versions of TCP behave the same when packet loss is detected by a retransmission timeout; TCP reduces its window to 1 packet and goes to slow-start phase until it reaches half the previous window, where it starts the congestion-avoidance phase. Some versions of TCP (e.g. TCP Reno, New-Reno [19]) attempt to prolong their operation in congestion-avoidance, but eventually go back to slow-start on a retransmission timeout.

Ideally, TCP would operate in the congestion-avoidance stage to efficiently utilize the network resources. From the

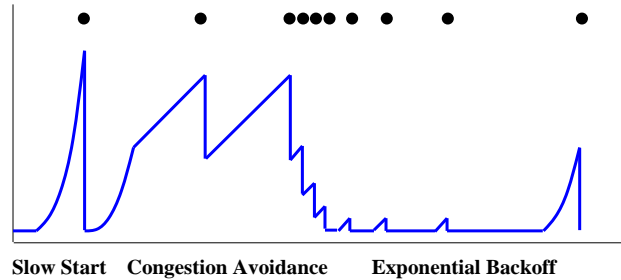


Fig. 1. Behavior of TCP congestion control

traffic characteristics point-of-view, TCP generates more stable traffic in this stage than the other two stages (slow-start and exponential-backoff). Unfortunately, two factors drag TCP away from the congestion-avoidance stage: (1) Internet measurements [20] show that most flows (both TCP and UDP) are short in size (less than 100 packets). As mentioned earlier, TCP's adaptive control mechanism requires a certain period of time to learn the state of the network, such as the round-trip time and average share of the available bandwidth. Because most of the TCP connections are too short to generate enough packet samples to obtain such detailed knowledge, most of their packets are sent out at a conservatively estimated rate during slow-start; and (2) an even worse situation is that some parts of the current Internet can be highly congested [21, 22], and thus the packet loss rate is relatively high. High loss rate makes consecutive packet losses possible and retransmission timeouts force many TCP connections to stay in the exponential-backoff stage.

Under these conditions, a single TCP connection may spend most or all of its time in slow-start or exponential-backoff stages. This is the situation we explore in the next section.

III. MODELING TCP BEHAVIOR IN SLOW-START AND EXPONENTIAL-BACKOFF STAGES

In this section, we use a discrete-time Markov chain to describe the behavior of TCP in slow-start and exponential-backoff stages and explain why it generates heavy tails in packet inter-arrival times. To this end, we assume that TCP's window adaptation policy depends only on the last value of the congestion window. We can then draw a Markov chain as shown in Figure 2. We define the state of the Markov chain to be the negation of the binary logarithm of the packet inter-departure time, normalized by the average round-trip time, for each outgoing packet. For example, the -1 state means the current packet is going to be sent out 2 round-trip times after the previous packet. Notice that for slow start phase (state index $k > 0$, there are more than one (up to the current window size W)

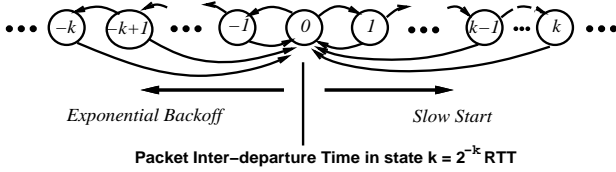


Fig. 2. Modeling the behavior of TCP congestion control

packet samples for each state. To simplify our analysis, we assume that all packets in the window are emitted evenly over the round trip time. Thus each state actually contains W mini-states, each of which corresponds to one packet in the current window, with inter-departure time equals to $1/W$ round trip time. However, for simplicity, we ignore the transition between these mini-states, but rather analyze the aggregated state with state index $\log_2 W$.

A. Modeling Slow-Start

Figure 3 illustrates the Markov chain for TCP slow-start stage. During this stage, TCP increases its window size

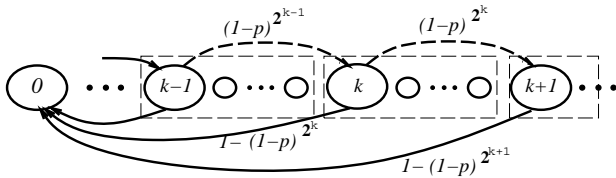


Fig. 3. Markov chain for slow-start

by one every time it gets an acknowledgment. Therefore essentially TCP doubles its window size every round-trip time if none of the packets in the previous window gets lost during the previous round-trip time. Assuming loss event is independent of the TCP behavior and the loss rate is constant p , then the probability of not losing any packet is $(1-p)^W$, where W is the previous round's window size. Recall that we define state k to denote a TCP window size of 2^k , or equivalently a rate of 2^k packets per RTT. Thus, the probability of transition from state k to $(k+1)$ equals to $(1-p)^{2^k}$, as shown in Figure 3. On the other hand, when packet loss is encountered, TCP reduces its window size to 1, and goes to slow-start again¹. Thus, the transition probability from state k to state 0 is $1 - (1-p)^{2^k}$.

Denote by π_k the probability that TCP is in state k , then we can write the following balance equations of the Markov chain:

$$\pi_k = (1-p)^{2^{k-1}} \pi_{k-1} \quad \text{for } k > 0$$

Therefore, we have

$$\pi_k = (1-p)^{2^{k-1}} \pi_{k-1}$$

¹This is how TCP Tahoe always behaves upon detecting a packet loss.

$$\begin{aligned} &= \dots \\ &= (1-p)^{2^{k(k-1)/2}} \pi_0 \quad \text{for } k > 0 \end{aligned} \quad (1)$$

B. Modeling Exponential-Backoff

The other part of the Markov chain is shown in detail in Figure 4. In this stage, TCP tries to retransmit the packet

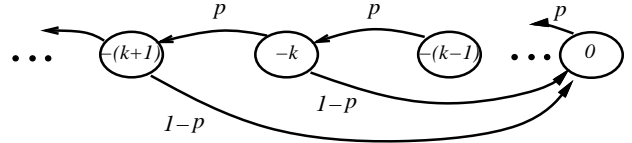


Fig. 4. Markov chain for exponential-timeout

lost in the previous round. The window value is kept at 1, and TCP employs Karn's algorithm to estimate the round-trip time, or the maximum time it needs to wait for the acknowledgment of the retransmitted packet. Instead of using the sampled round-trip time, TCP backs off the retransmission timeout (RTO) exponentially until the packet is acknowledged, usually starting from the latest RTT sample for a new (not-retransmitted) packet. If no RTT samples have yet been received, RTO is usually set to an initial default value.

During exponential-backoff, if TCP does not receive the acknowledgment of the retransmitted packet before RTO expires, either because the packet or its acknowledgment is lost (with probability p), or because the actual RTT is greater than RTO, then TCP doubles RTO and repeats the above process. The transition probability from state $-(k-1)$ to $-k$ in Figure 4 corresponds to this case.²

We can derive the state probabilities from the following balance equations:

$$\pi_{-k} = p\pi_{-(k-1)} \quad \text{for } k > 0$$

Therefore, we have:

$$\begin{aligned} \pi_{-k} &= p\pi_{-(k-1)} \\ &= \dots \\ &= p^k \pi_0 \quad \text{for } k > 0 \end{aligned} \quad (2)$$

C. Comments and Observations on the Analysis

Note that in Figure 2, we don't include the congestion-avoidance stage. Our objective here is to model *short* TCP connections. Many Internet measurements [5, 20] have shown that most of the Internet flows are short-lived. Since

²Note that in Karn's algorithm, if a packet has been sent more than once, the initial RTO for the next packet is set to the RTO of the previous transmission, not the estimated round-trip time. In other words, TCP stays in state $-k$ instead of going back to state 0. We ignore this detail to keep our Markovian analysis simple.

TCP needs to conservatively adapt to the available bandwidth in the network, it is often the case that short TCP connections end their transmission before they really enter into the congestion-avoidance stage. Current versions of TCP try to avoid being too pessimistic when a packet loss is encountered. They employ *fast retransmit and recovery* [19] mechanisms to avoid going to exponential-backoff stage too often. To do this, if the TCP sender receives the K^{th} duplicate acknowledgment packet before the timer expires,³ the sender will behave as if a timeout has occurred and retransmits while remaining in congestion-avoidance. However, for short connections, since they rarely have a chance to open their window big enough to have the fast retransmit mechanism workable, most of their lifetime is spent in exponential-backoff and slow-start. Thus, the model in Figure 2 is more suited for short TCP connections.

The steady-state probability distribution function for the states in the Markov chain is plotted in Figure 5:

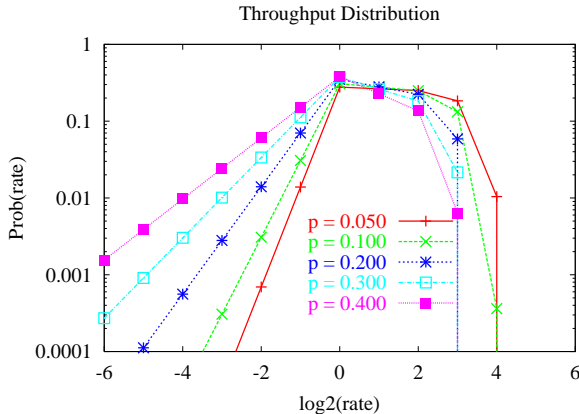


Fig. 5. State probability with different loss rate

Observing the part of the curve for the exponential-backoff stage (the left side), when p is large enough (e.g. $p > 12\%$), the distribution of packet inter-arrival times from the single TCP source starts to exhibit a tail consistent with infinite variance (if it is extended to infinity). This is shown in the following analysis. We have noticed from equation (2) that $\pi_{-k} = p^k \pi_0$, where π_{-k} corresponds to the probability of inter-arrival of (retransmitted) packets equals 2^k RTT. Denote by T the random variable which describes the packet inter-arrival time in the exponential-backoff stage, then we have:

$$\begin{aligned} Pr[T = t = 2^k] &= \pi_{-k} \\ &= p^k \pi_0 \\ &= (2^k)^{-\log_2(1/p)} \pi_0 \end{aligned}$$

³ K is usually 3.

$$\begin{aligned} &= t^{-\log_2(1/p)} \pi_0 \\ &= t^{-(\alpha+1)} \pi_0 \end{aligned} \quad (3)$$

where α is the slope (exponent) of the corresponding complementary cumulative distribution. Also, observe from Figure 5 that the window (burst) size distribution in the slow-start stage is *not* heavy-tailed distributed as the tail of the distribution decays very fast.

Willinger *et al.* [4] have shown that the result of aggregating many i.i.d. ON/OFF processes, where either the ON or OFF period of each process follows a heavy-tailed distribution, is a self-similar fractional Gaussian noise process, with $H = (3 - \min(\alpha_1, \alpha_2))/2$, where α_1 and α_2 are the shape parameters of the heavy-tailed distribution for ON and OFF periods respectively, given that $1 < \alpha_1, \alpha_2 < 2$. From equation (3), we can compute the shape parameter α as:

$$\alpha = \log_2\left(\frac{1}{p}\right) - 1 = \log_2\left(\frac{1}{2p}\right) \quad (4)$$

Therefore, we can compute the Hurst parameter of the time-series generated by our Markovian modeled processes as:

$$H = (3 - \alpha)/2 = (3 - \log_2\left(\frac{1}{2p}\right))/2 \quad (5)$$

when $1 < \alpha < 2$, or $\frac{1}{8} < p < \frac{1}{4}$ corresponding to $0.5 < H < 1$.

Notice that when p increases, the value of the Hurst parameter, or the degree of self-similarity, increases as well. This conclusion agrees with the actual Ethernet trace shown in [3]. Finally, as described above, note that the exponential backoff algorithm operates between a lower limit (the current retransmission timeout) and an upper limit. The upper limit is generally 64 times the lower limit. When combined with other sources of variability, this means that the scaling of packet inter-arrivals due to this effect occurs over a range of scales typically spanning about 3 decimal orders of magnitude.

IV. SIMULATION WITH PATHOLOGICAL TCP CONNECTIONS

A. Demonstrate Self-similarity by Simulation

We have set up a simple experiment with the NS [10] simulator to validate the above observation. As in many other similar studies [7, 22], we simulate the case that TCP connections go through a bottleneck link. However, because we are interested in the behavior of a single TCP flow, instead of simulating the case of many TCP flows competing for the bottleneck resource, we simulate only a

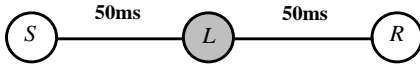


Fig. 6. Simulation topology

single connection transmitting over a lossy link. The topology is shown in Figure 6. Node S is the sender, node R is the receiver, the two-way propagation delay between S and R is set to 200 msec. We put a loss module inside node L , which drops incoming (new or retransmitted) packets uniformly with probability p . The packet size is set to be fixed at 576 bytes. The link capacity, buffer size and receiver window size are set to very large values so that only the loss module would affect TCP’s performance. The size of each TCP connection is set to be fixed at $FlowSize$ packets. A terminated connection is immediately replaced by a new connection (with all TCP parameters reset). We use Tahoe version of TCP.⁴

We have traced the packet arrival event during the simulation both at the upstream link (S, L) and downstream link (L, R) of node L .⁵ Figure 7 shows four time-series plots of the traffic on the upstream link. The plots are produced by aggregating packet traffic into discrete bins of 0.2, 2, 20, or 200 seconds. We observe that traffic is bursty over all time scales shown, and that the variability does not decline rapidly as time scale increases. This is suggestive of non-trivial self-similarity in the traffic trace.

We use the Aggregated Variance method as well as the Wavelet Analysis method to test the self-similarity of the 55-hour packet trace for the simulation with parameters $p = 0.15$ and $FlowSize = 15$ packets.⁶ Figure 8 shows the results. Roughly, with the Aggregated Variance method, for a self-similar process, the variance of the time average goes to zero very slowly as the bin size, or aggregation level, becomes larger. In Figure 8, the sample variance is plotted versus bin size on a log-log plot. For a self-similar process, the result should be a straight line

⁴We also simulated TCP Reno, the most widely used version of TCP. The simulation results are very similar to TCP Tahoe. Note that in the regime of high and bursty packet losses we are considering in this paper, TCP Tahoe is actually more robust to multiple losses than TCP Reno.

⁵All the data shown in this paper are collected at the upstream link (S, L), i.e., the data includes all packets, including those that are later lost at node L . Excluding those (eventually lost) packets affects the distributions (packet inter-arrival times, burst sizes) at small time scales, but does not affect the tail shape of the distributions (i.e., at large time scales). Thus the major conclusions of this paper still hold for packet traces collected on the downstream link (L, R), for which results are not shown.

⁶For different simulation setups, the packet trace contains hundreds of thousands to few millions of packets transmitted over the upstream link.

with a slope of $2H - 2$. If the time-series has no long-range dependence, then $H = 0.5$ and the slope should be -1 [23]. With the Wavelet Analysis Method, for a self-similar process, the degree of variability, or burstiness, decreases very slowly at larger time scales. In Figure 8, this variability is represented by “energy” estimates with 95% quantiles, computed using the tool of Veitch and Abry [24]. The plot of energy versus time-scale should be a straight line with a slope of $2H - 1$.⁷ The results of Figure 8 show values of H that are significantly higher than 0.5, which is again indicative of self-similarity.

B. Heavy Tails in Packet Inter-arrivals: How and Why

Greiner *et al.* show in [26] that a renewal process where the interarrival times have a power-tail yields self-similar scaling behavior. Our simulation results, which closely match our analysis in Section III, show that packet inter-arrivals under TCP can indeed be heavy-tailed, causing self-similarity in the overall traffic. Figure 9 shows simulation results for $p = 0.15$ and $FlowSize = 15$ packets.

We consider packets whose inter-arrival time is less than 2 times the packet transmission time (about 1 millisecond in our simulation) to constitute a burst. Bursts are only emitted in slow-start, so their presence does not affect the heavy tails in our Markov chain model. The probability of a state π_k ($k = 1, 2, 3$) can be thought of as the probability of a burst of size 2^k in our NS simulation. Therefore, we use two plots to represent the inter-packet arrival behavior. Figure 9(a) shows the log-log complementary distribution (LLCD) plot for those inter-arrival times bigger than 1 millisecond, and Figure 9(b) shows the distribution of packet arrival in bursts. Roughly, Figure 9(a) describes TCP behavior in the exponential-timeout stage (and part of the slow-start stage), while Figure 9(b) describes the slow-start stage. Note that Figures 9(a) and (b) are plotted in \log_{10} - \log_{10} and \log_{10} - \log_2 scales, respectively. Also note that each burst corresponds to an aggregated state in the Markov chain (with state index $k > 0$).

We observe that the inter-arrival time distribution in Figure 9(a) clearly has a tail that shows hyperbolic shape. The function given by the linear least-squares fit has the shape $\alpha = \log_2 3.33 = 1.74$ (implies high variance). The shape of the tail matches the value obtained by analysis from equation (4) in Section III-C.

The burst size distribution in the slow-start plot in Figure 9(b) also shows a shape similar to our analysis given in Figure 5, where it decreases very fast due to the relatively high loss probability.

⁷Note in Figure 8(b), the “dip” in the Wavelet plot reflects the strong correlation introduced by the network round-trip time (of 200 msec. here), as explained in [25].

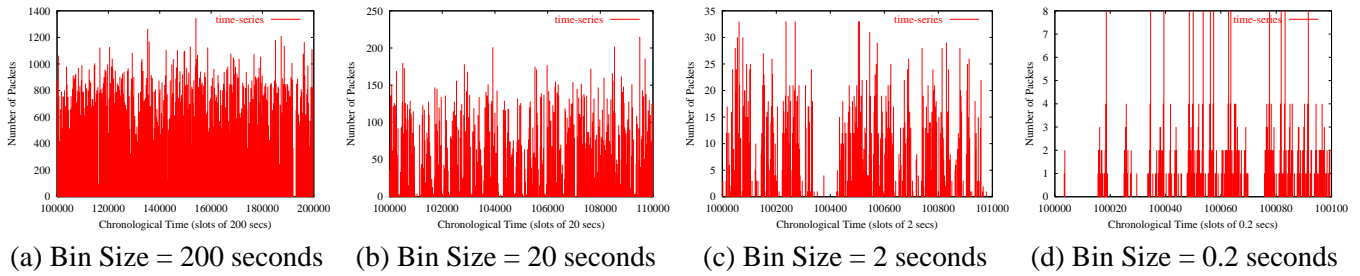
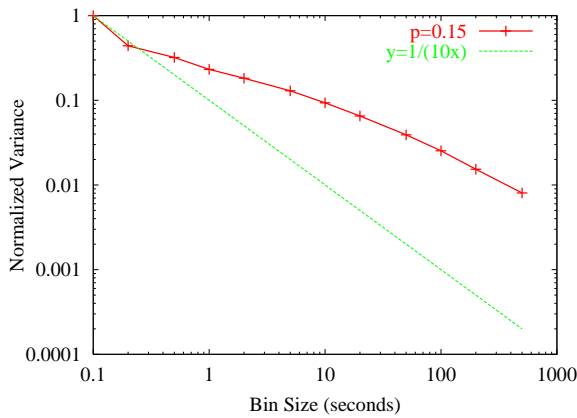
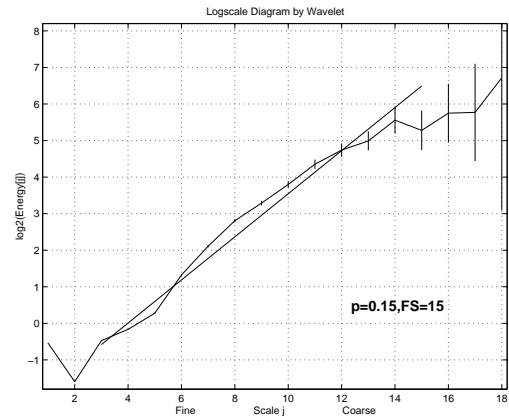


Fig. 7. Traffic burstiness over four orders of magnitude

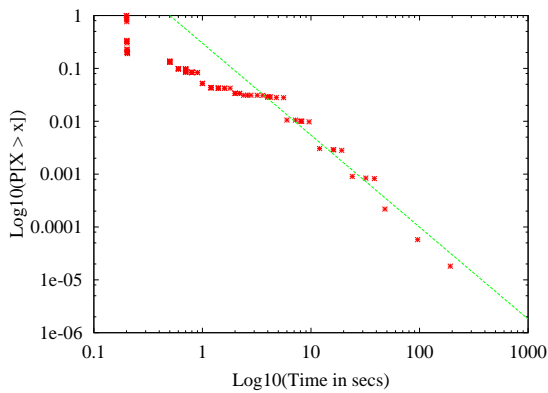


(a) Aggregated Variance Method $H = 0.76$

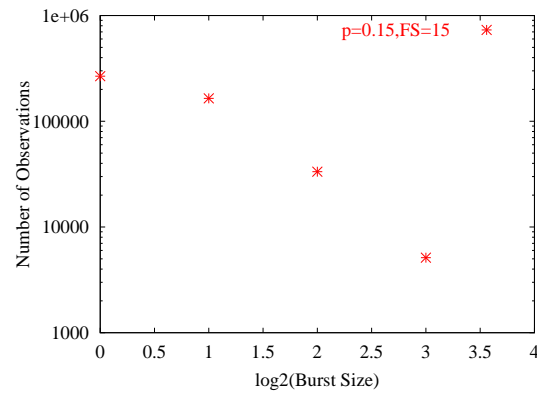


(b) Wavelet Method $H = 0.795 (0.793, 0.797)$

Fig. 8. LRD analysis on packet trace with loss rate = 15%



(a)



(b)

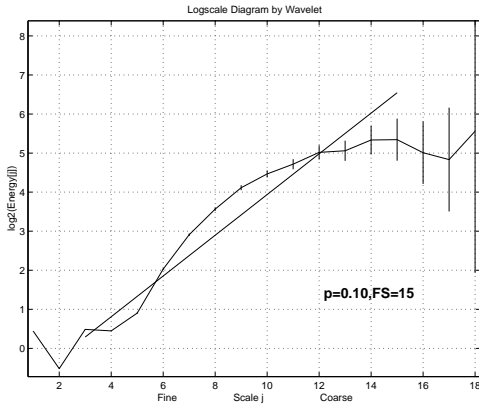
Fig. 9. The inter-arrival time and burst size distribution of TCP in (a) exponential-timeout (CDF), (b) slow-start stage (pdf)

C. Effect of Loss Rate

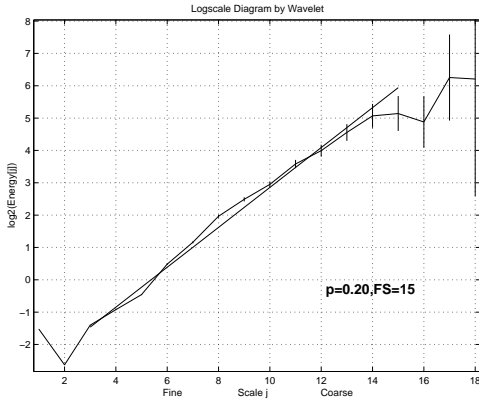
We also vary the loss rate p in our simulation experiments. Figure 10 plots the Wavelet analysis and the inter-arrival time distribution for $p = 0.1$ and $p = 0.2$. Observe that the slopes of the linear least-squares fit obtained by simulation closely match those obtained by analysis (cf. equation (4)).

When the loss rate is relatively low, the self-similarity exhibited by the time-series is less significant. This is because the tail of the inter-arrival time (or the OFF period) distribution is less heavy and thus at larger time scales,

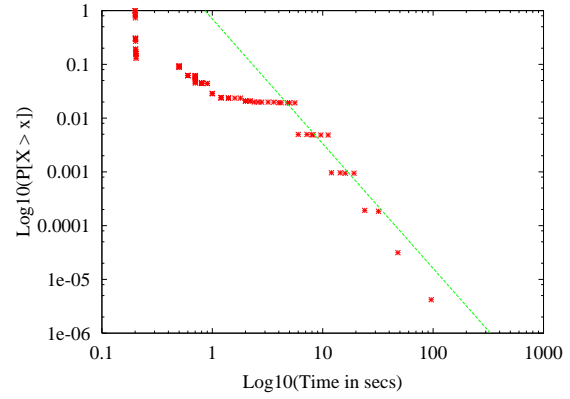
the traffic starts to show less variance. On the contrary, with higher loss rate, the tail of the OFF periods is heavier and the statistical self-similarity is more significant. Recall that in our analysis (cf. equation (5)), we require $p > 1/8$ so that the tail of the OFF period distribution is heavy enough to introduce self-similarity in the overall traffic (i.e. $H > 0.5$). However, in our simulation, we observe reasonably high variability at a lower p of 0.1. That indicates that our analysis may underestimate the degree of self-similarity or long-range dependence that may be present. The reason is due to the variability in the initial



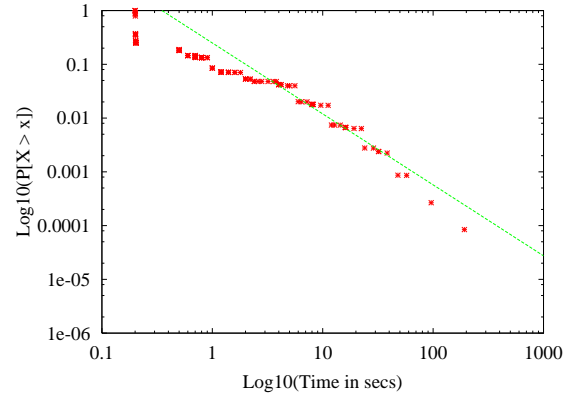
(a) $p = 0.1$, Wavelet Method $H = 0.761$



(c) $p = 0.2$, Wavelet Method $H = 0.808$



(b) LLCD of inter-arrival times, slope = $-\log_2 5 = -2.32$.



(d) LLCD of inter-arrival times, slope = $-\log_2 2.5 = -1.32$.

Fig. 10. Simulation with loss rates $p = 0.1$ (a,b), and $p = 0.2$ (c,d)

timeout values for different packets that we observe in our simulation, and that we do not capture in our analytical model. This variability causes the traffic to have a high variance at large time scales even for values of p smaller than $1/8$. Section IV-D discusses in more detail the effect of variability in initial timeout values.

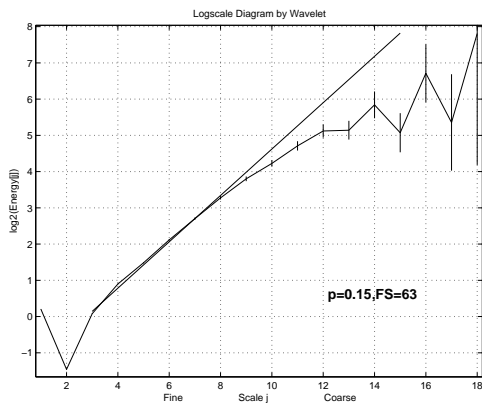
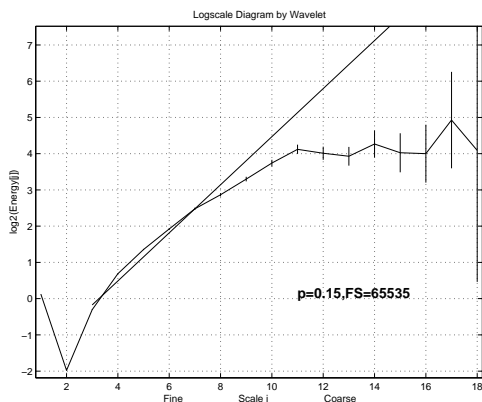
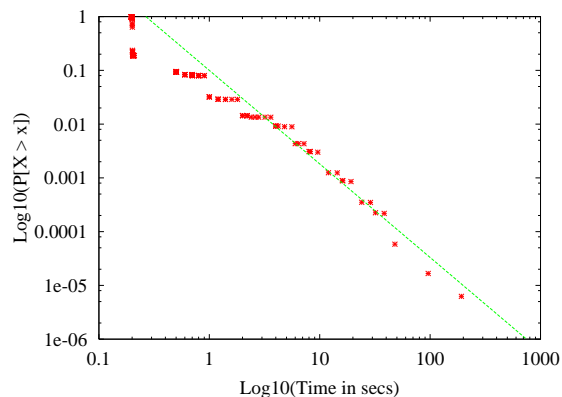
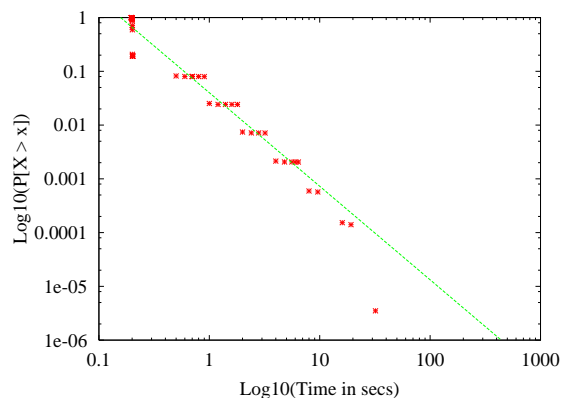
Our analysis in Section III-C estimates the Hurst parameter H from the shape parameter α (cf. equation (5)). However, the Wavelet method shows a different Hurst value for some values of loss rate. For example, when $p = 0.15$ and $\alpha = 1.74$, our analysis yields $H = (3 - 1.74)/2 = 0.63$, whereas the Wavelet method yields $H = 0.795$. However, when the loss rate is high, our analysis gives a better estimation. For example, when $p = 0.2$ and $\alpha = 1.32$, we obtain $H = 0.84$, which is very close to $H = 0.808$ obtained using the Wavelet method. We conjecture that when the loss rate is low, the effect of variable initial retransmission timeout values for different packets is more significant in determining the Hurst parameter.

D. The Variability of Initial Retransmission Timeout Value

There are two reasons to have variable initial timeout. The first is brought by variable queuing delays. TCP in-

cludes the variance of the RTT in its RTT estimate, which is then used to compute the retransmission timeout value. The second factor, which is more significant in our model, is that for the first packet, the initial timeout value is set to a default value because TCP does not have any packet sample to estimate the round-trip time. This default value is suggested to be 3 seconds in [27], and is set to 6 seconds by default in the NS simulator. Recall that the two-way propagation delay in our simulation is 0.2 second, so if TCP had enough packet samples to accurately estimate the round-trip time, the initial retransmission timeout value would have been set to a value close to 0.4 second.⁸ This huge difference (about 15 times) results in some packet inter-arrival times to be as high as $64 \times 6 = 384$ seconds, if the transmission of a packet fails 6 times. The effect of the default initial timeout value is less significant in two cases: when its value is closer to the estimated value, and when the connection size is large. For these two cases, we conducted the following two experiments.

⁸Observe that Figure 9 shows a minimum packet inter-arrival time of 0.2 second.

(a) Wavelet Method $H = 0.820$ (c) Wavelet Method $H = 0.832$ (b) LLCD of inter-arrival times, slope = $-\log_2 3.33 = -1.74$ (d) LLCD of inter-arrival times, slope = $-\log_2 3.33 = -1.74$ Fig. 11. Simulation with file sizes $FlowSize = 63$ (a,b), and $FlowSize = 65535$ (c,d)

D.1 Effect of File (Flow) Size

We compare the packet traces generated from sending different file sizes over TCP. We fix the loss rate at 15%. Figure 11 shows both the Wavelet analysis on the packet trace and the packet inter-arrival time distribution for $FlowSize = 63$ packets and 65, 535 packets.

We observe that as the file size increases, the time-series shows less variations at larger time scales; this is apparent from the less heavy tail of the LLCD plot and the Wavelet plot bending down at large time scales. Note that the high value of $H = 0.832$ in Figure 11(c) is misleading as the linear least-squares fit does not match quite well the energy curve. Indeed, the time-series produced by transmitting large files is no longer self-similar. We attribute this reduced variability to the more accurate information TCP learned from increased packet samples and better rate estimation (through congestion-avoidance).

D.2 Effect of Default Initial Timeout Value

We conducted two experiments to reduce the difference between the default initial timeout value and the normal timeout value. In the first experiment, we increase the two-

way propagation delay from 0.2 second to 1 second. In the second experiment, we reduce the default initial timeout value from 6 seconds to 0.5 second. Figure 12 shows the Wavelet analysis of the time-series produced by the captured packet traces for $FlowSize = 15$ packets and $p = 0.15$.

As the variance of the timeout values becomes smaller, the long-range dependence in the time-series becomes less significant.⁹ In [7], the authors show that when 40 TCP connections are squeezed into a small pipe (on average, each connection can transmit less than one packet in one round-trip time), the time-series produced by each connection shows long-range dependence. We can explain this phenomena here based on the above observations. First of all, a highly congested network brings higher packet losses (indeed, [7] considers a loss rate of around 16%). Secondly, in this case, queueing delays bring more variability to the initial retransmission timeout values for different packets. We observe that it is these two factors that contribute to heavy-tailed packet inter-arrival times, or the

⁹This observation agrees with [16], in which Peha observes that for larger or more random timeouts, self-similar behavior is observed at larger time scales.

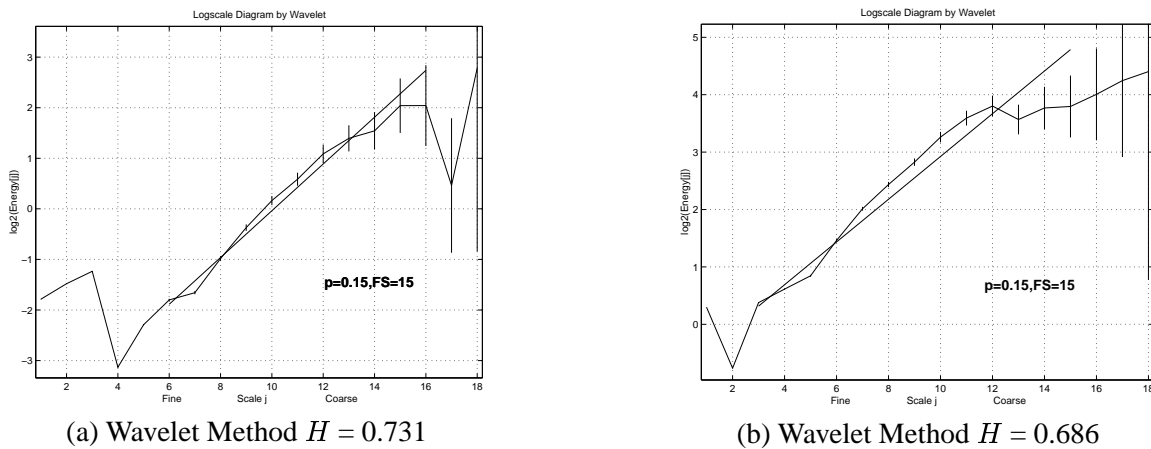


Fig. 12. Reducing the difference between timeout values: (a) increase round-trip propagation delay to 1 second, (b) decrease default timeout value to 0.5 second.

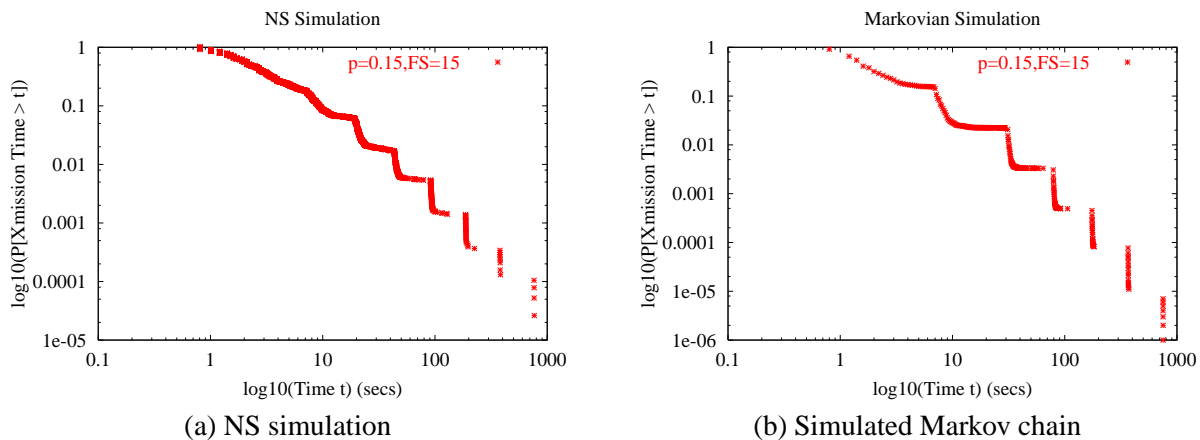


Fig. 13. Validation of heavy-tailed transmission time, $FlowSize = 15$, $p = 15\%$

“chaotic” nature of the TCP congestion control.

E. Network Dynamics and Heavy-tailed Transmission Times

In [8], the authors noted that the omnipresence of heavy-tails in file size distribution in the majority of file systems is not the only reason behind the heavy-tailed distribution of HTTP file transfers. They pointed out that the network condition is also important by observing that there is *no* strong correlation between file transmission time and file size. However, that paper made no attempt to explore alternate reasons for heavy-tailed transmission times.

We observe that in addition to these reasons, TCP’s adaptive congestion control mechanism in the presence of severe network conditions can be another major contributor to heavy tails. We run simulations on both the NS simulator and the Markov chain described in Section II to validate this observation. Figure 13 shows the LLCD plot of the transmission time distribution for transferring a 15-packet file over a lossy channel with 15% loss rate.

The stair-step shape of the plot is due to the small file size (15 packets) leading to a small set of possible packet-loss events and thus transmission times. Nonetheless it shows an approximate power-law shape (the general trend apart from the stair-stepping) that is indicative of heavy tails in transmission times. Larger files would presumably show a smoother curve. This surprising effect means that under certain network conditions, heavy-tailed transmission times can occur even in the absence of *any* variability in file sizes.¹⁰

V. CONCLUSION AND FUTURE WORK

The principle of TCP congestion control is to avoid overloading the network while still maintaining network usage efficiency. The TCP control algorithm is simple but effective. However, this simplicity does not come for free. The behavior of TCP becomes less predictable or even

¹⁰This high variability in file transmission times is consistent with that observed in [6].

chaotic when the network condition is out of TCP's control. In this paper, we have demonstrated that when a TCP connection is going through a highly lossy channel and the loss condition is not affected by this single TCP connection's behavior, TCP starts to produce packet trains that show self-similarity. We used a simple Markovian model to demonstrate why heavy-tails show up in packet inter-arrival times and file transmission times. We believe that our analysis sheds some light on the relationship between TCP conservative control mechanism and heavy-tailed distributions observed in Internet traffic.

Our future work includes extending our analytical model to consider the variability in round-trip time estimation, and verifying our analytical model in more complex network settings. We also plan to study possible improvements to TCP congestion control mechanism under severe network conditions. We are also planning to investigate the relationship between network-level self-similarity and application-level self-similarity (especially for network-aware adaptive applications).

REFERENCES

- [1] M. Crovella and A. Bestavros, "Self-Similarity in World Wide Web Traffic: Evidence and Possible Causes," *IEEE/ACM Transactions on Networking*, pp. 835–846, December 1997.
- [2] A. Feldmann, A.C. Gilbert, W. Willinger, and T.G. Kurtz, "The Changing Nature of Network Traffic: Scaling Phenomena," *ACM Computer Communication Review*, pp. 5–29, April 1998.
- [3] W.E. Leland, M.S. Taqqu, W. Willinger, and D.V. Wilson, "On the Self-Similar Nature of Ethernet Traffic (Extended Version)," *IEEE/ACM Transactions on Networking*, pp. 1–15, February 1994.
- [4] W. Willinger, V. Paxson, and M.S. Taqqu, "Self-Similarity and Heavy-Tails: Structural Modeling of Network Traffic," in *A Practical Guide To Heavy Tails: Statistical Techniques and Applications*, R.J. Adler, R.E. Feldman and M.S. Taqqu, editors. ISBN 0-8176-3951-9. Birkhäuser, Boston, 1998.
- [5] K. Thompson, G.J. Miller, and R. Wilder, "Wide-Area Internet Traffic Patterns and Characteristics," *IEEE/ACM Transactions on Networking*, pp. 10–23, November 1997.
- [6] R.T. Morris, *Scalable TCP Congestion Control*. Ph.D. thesis, Harvard University, Cambridge MA, The Division of Engineering and Applied Sciences, January 1999.
- [7] A. Veres, B. Hungary, and M. Boda, "The Chaotic Nature of TCP Congestion Control," in *Proceedings of IEEE INFOCOM'2000*, March 2000.
- [8] M.E. Crovella, M.S. Taqqu, and A. Bestavros, "Heavy-Tailed Probability Distributions in the World Wide Web," in *A Practical Guide To Heavy Tails: Statistical Techniques and Applications*, R.J. Adler, R.E. Feldman and M.S. Taqqu, editors. ISBN 0-8176-3951-9. Birkhäuser, Boston, 1998.
- [9] A. Feldmann, *On-line Call Admission for High-Speed Networks*, Ph.D. thesis, Carnegie Mellon University, Pittsburgh, PA, School of Computer Science, September 1995.
- [10] UCB, LBNL, and VINT, "Network simulator - ns (version 2)," <http://www-mash.cs.berkeley.edu/ns/>, 1999.
- [11] R.E. Feldman R.J. Adler and M.S. Taqqu, *A Practical Guide To Heavy Tails: Statistical Techniques and Applications*, ISBN 0-8176-3951-9. Birkhäuser, Boston, 1998.
- [12] Jan Beran, "Statistics for Long-Memory Processes," in *Mono-graphs on Statistics and Applied Probability*, Chapman and Hall, New York, NY, 1994.
- [13] W. Willinger, M.S. Taqqu, R. Sherman, and D.V. Wilson, "Self-Similarity Through High-Variability: Statistical Analysis of Ethernet LAN Traffic at the Source Level," *IEEE/ACM Transactions on Networking*, pp. 71–86, February 1997.
- [14] V. Paxson and S. Floyd, "Wide Area Traffic: The Failure of Poisson Modeling," *IEEE/ACM Transactions on Networking*, pp. 236–244, June 1995.
- [15] J.H.B. Deane, C. Smythe, and D.J. Jefferies, "Self-Similarity in a Deterministic Model of Data Transfer," *Journal of Electronics*, vol. 80, no. 5, pp. 677–691, 1996.
- [16] J. M. Peha, "Retransmission Mechanisms and Self-Similar Traffic Models," in *Proceedings of IEEE/ACM/SCS Communication Networks and Distributed Systems Modeling and Simulation Conference*, Phoenix, Arizona, Jan. 1997, pp. 47–52.
- [17] A. Veres, Zs. Kenesi, S. Molnár, and G. Vattay, "On the Propagation of Long-Range Dependence in the Internet," in *Proceedings of ACM SIGCOMM 2000*, Stockholm, Sweden, Aug.-Sep. 2000.
- [18] P. Karn and C. Partridge, "Improving Round-Trip Time Estimates in Reliable Transport Protocols," *ACM Transactions on Computer Systems (TOCS)*, vol. 9, pp. 365–373, 1991.
- [19] K. Fall and S. Floyd, "Simulation-based Comparisons of Tahoe, Reno and SACK TCP," *ACM Computer Communication Review*, pp. 5–21, July 1996.
- [20] K. Claffy, G. Miller, and K. Thompson, "the Nature of the Beast: Recent Traffic Measurements from an Internet Backbone," in *Proceedings of INET'98*, <http://www.caida.org/outreach/papers/Inet98> 1998.
- [21] V. Paxson, "End-to-End Internet Packet Dynamics," in *Proceedings of ACM SIGCOMM'97*, Cannes, France, September 1997.
- [22] R.T. Morris, "TCP Behavior with Many Flows," in *IEEE International Conference on Network Protocols (ICNP'97)*, Atlanta, Georgia, October 1997.
- [23] Murad Taqqu, "Statistical Methods for Long-Range Dependence," <http://math.bu.edu/people/murad/methods/index.html>.
- [24] D. Veitch and P. Abry, "A Wavelet Based Joint Estimator of the Parameters of Long-Range Dependence," *IEEE Transactions on Information Theory*, April 1999.
- [25] A. Feldmann, A.C. Gilbert, P. Huang, and W. Willinger, "Dynamics of IP traffic: A Study of the Role of Variability and the Impact of Control," in *Proceeding of ACM/SIGCOMM'99*, Boston, Mass, September 1999.
- [26] M. Greiner, M. Jobmann, and L. Lipsky, "The Importance of Power-Tail Distributions for Modeling Queueing Systems," *Operations Research*, vol. 47, no. 2, March-April 1999.
- [27] R. Braden, "Requirements for Internet Hosts – Communication Layers," Internet RFC 1122, October 1989.