# Internet Traffic Tends *Toward* Poisson and Independent as the Load Increases

Jin Cao, William S. Cleveland, Dong Lin, Don X. Sun

*Abstract*— **Network devices put packets on an Internet link, and multiplex, or superpose, the packets from different active connections.**

**Extensive empirical and theoretical studies of packet traffic variables — arrivals, sizes, and packet counts — demonstrate that the number of active connections has a dramatic effect on traffic characteristics. At low connection loads on an uncongested link — that is, with little or no queueing on the link-input router — the traffic variables are long-range dependent, creating burstiness: large variation in the traffic bit rate. As the load increases, the laws of superposition of marked point processes push the arrivals toward Poisson, the sizes toward independence, and reduces the variability of the counts relative to the mean. This begins a reduction in the burstiness; in network parlance, there are multiplexing gains.**

**Once the connection load is sufficiently large, the network begins pushing back on the attraction to Poisson and independence by causing queueing on the link-input router. But if the link speed is high enough, the traffic can get quite close to Poisson and independence before the push-back begins in force; while some of the statistical properties are changed in this high-speed case, the push-back does not resurrect the burstiness. These results reverse the commonly-held presumption that Internet traffic is everywhere bursty and that multiplexing gains do not occur.**

**Very simple statistical time series models — fractional sum-difference (FSD) models — describe the statistical variability of the traffic variables and their change toward Poisson and independence before significant queueing sets in, and can be used to generate open-loop packet arrivals and sizes for simulation studies.**

**Both science and engineering are affected. The magnitude of multiplexing needs to become part of the fundamental scientific framework that guides the study of Internet traffic. The engineering of Internet devices and Internet networks needs to reflect the multiplexing gains.**

## I. ARE THERE MULTIPLEXING GAINS?

When two hosts communicate over the Internet — for example, when a PC and a Web server communicate for the purpose of sending a Web page from the server to the PC — the two hosts set up a *connection*. One or more files are broken up into pieces, headers are added to the pieces to form packets, and the two

hosts send packets to one another across the Internet. When the transfer is completed, the connection ends.

The headers, typically 40 bytes in size, contain much information about the packets such as their source, destination, size, etc. In addition there are 40-byte control packets, all header and no file data, that transfer information form one host to the other about the state of the connection. The maximum amount of file information allowed in a packet is 1460 bytes, so packets vary in size from 40 bytes to 1500 bytes.

Each packet travels across the Internet on a path made up of devices and transmission links between these devices. The devices are the two hosts at the ends and routers in-between. Each device sends the packet across a transmission link to the next device on the path. The physical medium, or the "wire", for a link might be a telephone wire from a home to a telephone company, or a coaxial cable in a university building, or a piece of fiber connecting two devices on the network of an Internet service provider. So each link has two devices, the sending device that puts the bits of the packet on the link, and the receiving device, which receives the bits. Each router serves as a receiving device for one or more input links and as the sending device for one or more output links; it receives the packet on an input link, reads the header to determine the output link, and sends bits of the packet.

Each link has a speed: the rate at which the bits are put on the wire by the sending device and received by the receiving device. Units are typically kilobits/sec (kbps), megabits/sec (mbps), or gigabits/sec (gbps). Typical speeds are 56 kbps, 1.5 mbps, 10 mbps, 100 mbps, 156 mbps, 622 mbps, 1 gbps, and 2.5 gbps. The *transmission time* of a packet on a link is the time it takes to put all of the bits of the packet on the link. For example, the transmission time for a 1500 byte (12000 bit) packet is 120 $\mu$s at 100 mbps and 12 $\mu$s at 1 gbps. So packets pass more quickly on a higher-speed link than on a lower-speed one.

The packet traffic on a link can be modeled as a marked point process. The arrival times of the process are the arrival times of the packets on the link;

a packet arrives at the moment its first bit appears on the link. The marks of the process are the packet sizes. An Internet link typically carries the packets of many active connections between pairs of hosts. The packets of the different connections are intermingled on the link; for example, if there are three active connections, the arrival order of 10 consecutive packets by connection number might be 1, 1, 2, 3, 1, 1, 3, 3, 2, and 3. This intermingling is referred to as "statistical multiplexing" in the Internet engineering literature, and as "superposition" in the literature of point processes.

If a link's sending device cannot put a packet the link because it is busy with one or more other packets that arrived earlier, then the device puts the packet in a queue, physically, a buffer. Queueing on the device delays packets, and if it gets bad enough, and the buffer size is exceeded, packets are dropped. This reduces the quality of Internet applications such as Web page downloads and streaming video. Consider a specific link. Queueing of the packet in the buffer of the link's sending device is *upstream queueing*; so is queueing of the packet on sending devices that processed the packet earlier on its flight from sending host to receiving host. Queueing of the packet on the receiving device, as well as on devices further along on its path is *downstream* queueing.

All along the path from one host to another, the statistical characteristics of the packet arrivals and their sizes on each link affect the downstream queueing, particularly the queueing on the link receiving device. The most accommodating traffic would have arrivals and sizes on the link that result in a traffic rate in bits/sec that is constant; this would be achieved if the packet sizes were constant (which they are not) and if they arrived at equally spaced points in time (which they do not). In this case we would know exactly how to engineer a link of a certain speed; we would allow a traffic rate equal to the link speed. There would be no queueing and no device buffers. The utilization, the ratio of the traffic rate divided by the link speed, would be 100%, so the transmission resources would be used the most efficiently. If the link speed were 1.5 mbps, the traffic rate would be 1.5 mbps.

Suppose instead, that the traffic is stationary with Poisson arrivals and independent sizes. There would be queueing, so a buffer is needed. Here is how we would engineer the link to get good performance. Suppose the speed is 1.5 mbps. We choose a buffer size so that a packet arriving when the queue is nearly full would not have to wait more than about 500 ms; for 1.5 mbps this would be about 100 kilobytes, or

800 kilobits. An amount of traffic is allowed so that only a small percentage of packets are dropped, say 0.5%. For this Poisson and independent traffic, we could do this and and achieve a utilization of 95%, so the traffic rate would be 1.425 mbps.

Unfortunately, we do not get to choose the traffic characteristics. They are dictated by the engineering protocols that underlie the Internet. What can occur is far less accommodating than traffic that has a constant bit rate, or traffic that is Poisson and independent. The traffic can be very *bursty*. This means the following. The packet sizes and inter-arrival times are sequences that we can treat as time series. Both sequences can have persistent, long-range dependence; this means the autocorrelations are positive and fall of slowly with the lag $k$, for example, like $k^{-\alpha}$ where $0 < \alpha \leq 1$. Long-range dependent time series have long excursions above the mean and long excursions below the mean. Furthermore, for the sizes and inter-arrivals, the coefficient of variation, the standard deviation divided by the mean, can be large, so the excursions can be large in magnitude as well as long in time. The result is large downstream queue-height distributions with large variability. Now, when we engineer a link of 1.5 mbps, utilizations would be much lower, about 40%, which is a traffic rate of 0.6 mbps.

Before 2000, this long-range dependence had been established for links with relatively low link speeds and therefore low numbers of simultaneous active connections, or connection loads, and therefore low traffic rates. But beginning in 2000, studies were undertaken to determine if on links with higher speeds, and therefore greater connection loads, there were effects due to the increased statistical multiplexing. Suppose we start out with a small number of active connections. What happens to the statistical properties of the traffic as we increase the connection load? In other words, what is the effect of the increase in magnitude of the multiplexing? We would expect that the statistical properties change in profound ways, not just simply that the mean of the inter-arrivals decreases. Does the long-range dependence dissipate? Does the traffic tend toward Poisson and independent, as suggested by the superposition theory of marked point processes? This would mean that the link utilization resulting from the above engineering method increase. In network parlance, are there would be multiplexing gains.

In this article, we review the results of the new studies on the effect of increased statistical multiplexing on the statistical properties of packet traffic on an Internet link.

## II. THE VIEW OF THE INTERNET CIRCA 2000

The study of Internet traffic beginning in the early 1990s resulted in extremely important discoveries in two pioneering articles [1], [2]: counts of packet arrivals in equally-spaced consecutive intervals of time are long-range dependent and have a large coefficient of variation (ratio of the standard deviation to the mean), and packet inter-arrivals have a marginal distribution that has a longer tail than the exponential. This means the arrivals are not a Poisson process because the counts of a Poisson are independent and the inter-arrivals are exponential. The title of the second article, "Wide-Area Traffic: The Failure of Poisson Modeling", sent a strong message that the old Poisson models for voice telephone networks would not do for the emerging Internet network. And because queue-height distributions for long-range dependent traffic relative to the average bit/rate are much greater than for Poisson processes, it sent a signal that Internet technology would have to be quite different from telephone network technology. The discovery of long-range dependence was confirmed in many other studies (e.g., [3], [4], [5]). The work on long-range dependence drew heavily on the brilliant work of Mandelbrot [6], both for basic concepts and for methodology.

Models of source traffic were put forward to explain the traffic characteristics [3], [7], [8], [9]. The sizes of transferred files utilizing a link vary immensely; to a good approximation, the upper tail of the file size distribution is Pareto with a shape parameter that is often between 1 and 2, so the mean exists but not the variance. A link sees the transfer of files whose sizes vary by many orders of magnitude. Modeling the link traffic began with an assumption of a collection of on-off traffic sources, each on (with a value of 1) when the source was transferring a file over the link, and off (with a value of 0) when not. Since the model has no concept of packets, just connections, multiplexing becomes summation; the link traffic is a sum, or aggregate, of source processes. Because of the heavy tail of the on process, the summation is long-range dependent, and for a small number of source processes, has a large coefficient of variation. We will refer to this as the *on-off aggregation theory*.

Before 2000, there was little empirical study of packet arrivals and sizes. Most of the intuition, theory, and empirical study of the Internet was based on a study of packet and byte counts. It took some time for articles to appear in the literature showing packet inter-arrivals and packet sizes are long-range dependent, although one might have guessed this from the results for counts. The first report in the literature of which we are aware appeared in 1999 [5]. The first articles of which we are aware that sizes are long-range dependent appeared in 2001 [10], [11].

While there was no comprehensive empirical study of the effect of multiplexing, before 2000 there were theoretical investigations. Some of the early, foundations-setting articles on Internet traffic contained conjectures that multiplexing gains did not occur. Leland *et al.* [1] wrote:
We demonstrate that Ethernet LAN traffic is statistically *self-similar*, ... and that aggregating streams of such traffic typically intensifies the self-similarity ('burstiness') instead of smoothing it.
Crovella and Bestavros [3] wrote:
One of the most important aspects of self-similar traffic is that there is is no characteristic size of a traffic burst; as a result, the aggregation or superposition of many such sources does not result in a smoother traffic pattern.
Further consideration and discussion however suggested that issues other than long-range dependence needed to be considered. Erramilli *et al.* [12] wrote
... the FBM [fractional Brownian motion] model does predict significant multiplexing gains when a large number of independent sources are multiplexed, the relative magnitude is reduced by $\sqrt{n}$ ....
Floyd and Paxson [7] wrote:
... we must note that it remains an open question whether in highly aggregated situations, such as on Internet backbone links, the correlations [of long-range dependent traffic], while present, have little actual effect because the variance of the packet arrival process is quite small.
In addition, there were theoretical discussions of the implications of increased multiplexing on queueing [13], [14], [15], [16], [17]. But the problem with such theoretical study is that results depend on the assumptions about the individual traffic sources being superposed, and different plausible assumptions lead to different results. Without empirical study, it was not possible to resolve the uncertainty about assumptions.

With no clear empirical study to guide judgment, many subscribed to a presumption that multiplexing gains did not occur, or were too small to be relevant. For example, Listani *et al.* [18] wrote:
... traffic on Internet networks exhibits the same characteristics regardless of the number of simultaneous sessions on a given physical link.

Internet service providers acted on this presumption in designing and provisioning networks, and equipment designers acted on it in designing devices.

## III. FOUNDATIONS: THEORY AND EMPIRICAL STUDY

Starting in 2000, a current of research was begun to determine the effect of increased multiplexing on the statistical properties of many Internet traffic variables, to determine if multiplexing gains occurred [10], [19], [20], [21], [22].

The empirical study of byte and packet counts of previous work was enlarged to include a study of arrivals and sizes. Of course, much can be learned from counts, but arrivals and sizes are the more fundamental traffic variables. It is arriving packets with varying sizes that network devices process, not aggregations of packets in fixed intervals, and packet and byte counts are derived from arrivals and sizes, but not conversely.

In keeping with a focus on arrivals and sizes, the superposition theory of marked point processes became a guiding theoretical framework, replacing the on-off aggregation theory that was applicable to counts but not arrivals and sizes [23], [24]. The two theories are quite different. For the on-off aggregation theory, one considers a sum of independent random variables, and a central limit theorem shows the limit is a normal distribution. For the superposition theory, in addition to the behavior of sums, one considers a superposition of independent marked point processes, and a central limit theorem shows the limit is a Poisson point process with independent marks, and quite importantly, *the theorem applies even when the inter-arrivals and marks of each superposed source point process are long-range dependent*.

The following discussion draws largely on the very detailed account in [19]. We will consider packet arrivals and sizes, and packet counts in fixed intervals. We omit the discussion of byte counts since their behavior is much like that of the packet counts.

## IV. THEORY: POISSON AND INDEPENDENCE

Let $a_j$, for $j = 1, 2, \ldots$ be the arrival times of packets on an Internet link where $j = 1$ is for the first packet, $j = 2$ is for the second packet, and so forth. Let $t_j = a_{j+1} - a_j$ be the inter-arrival times, and let $q_j$ be the packet sizes. We treat $a_j$ and $q_j$ as a marked point process. $a_j$, $t_j$, and $q_j$ are studied as time series in $j$. Suppose we divide time into equally-spaced intervals, $[\Delta i, \Delta(i + 1))$, for i = 1, 2, $\ldots$ where $\Delta$ might be

1 ms or 10 ms or 100 ms. Let $p_i$ be the packet count, the number of arrivals in interval $i$. The $p_i$ are studied as a time series in $i$.

Suppose the packet traffic is the result of multiplexing $m$ traffic sources on the link. Each source has packet arrival times, packet sizes, and packet counts. The arrival times $a_j$ and the sizes $q_j$ of the superposition marked point process result from the superposing of the arrivals and sizes of the $m$ source marked point processes. The packet count $p_i$ of the superposition process in interval $i$ results from summing the $m$ packet counts for the $m$ sources in interval $i$; theoretical considerations for the $p_i$ are, of course, the same as those for the on-off aggregation theory described earlier.

*Provided certain assumptions hold*, the superposition theory of marked point processes prescribes certain behaviors for $a_j$, $t_j$, $q_j$, and $p_i$ as $m$ increases [24]. The arrivals $a_j$ tend toward Poisson, which means the inter-arrivals $t_j$ tend toward independent and their marginal distribution tends toward exponential. The sizes $q_j$ tend toward independent, but there is no change in their marginal distribution. As discussed earlier, the $t_j$ and $q_j$ have been shown to be long-range dependent for small $m$. Thus the theory predicts that the long-range dependence of the $t_j$ and the $q_j$ dissipates. But the autocorrelation of the packet counts $p_i$ does not change with $m$ so its long-range dependence is stable. However, the standard deviation relative to the mean, the coefficient of variation, falls off like $1/\sqrt{m}$. This means that the burstiness of the counts dissipates as well; the durations of excursions of $p_i$ above or below the mean, which are long because of the long-range dependence, do not change because the correlation stays the same, but the magnitudes of the excursions get smaller and smaller because the statistical variability decreases.

The following assumptions for the source packet processes lead to the above conclusions:
• homogeneity: they have the same statistical properties.
• stationarity: their statistical properties do not change through time.
• independence: they are independent of one another and the size process of each is independent of the arrival process.
• non-simultaneity: the probability of two or more packet arrivals for a source in an interval of length $w$ is $o(w)$ where $o(w)/w$ tends to zero as $w$ tends to zero.

We cannot take the source processes to be the individual connections; they are not stationary, but rather

transient, that is, that have a start time and a finish time. Instead, we randomly assign each connection, to one of $m$ source processes. Suppose the start times are a stationary point process, and let $\rho$ be the arrival rate. Then the arrival rate for each source process is $\rho/m$. We let $\rho \to \infty$, keeping $\rho/m$ fixed to a number sufficiently large that the source processes are stationary; so $m \to \infty$.

We refer to the formation of the source processes, the assumptions about them, and the implications, as the *superposition theory*. It is surely true that all we have done with this theory is to reduce our uncertainty about whether the superposition process is attracted to Poisson and independent with an uncertainty about whether the above construction creates source processes that satisfy the assumptions. But it is a least plausible, although by no means certain, that there are cases where the source process satisfies the above assumptions over a range of values of $m$. What we have done is to create a plausible hypothesis to be tested by empirical study which we describe shortly.

## V. THEORY: THE NETWORK PUSHES BACK

While we cannot verify the hypotheses of the superposition theory without empirical study, we can at least quite convincingly describe a way in which the network can push back and defeat assumptions. Once $m$ is large enough, significant link-input queueing begins, and then grows as $m$ gets larger still; at some point, the queueing will be large enough that the assumptions of independence of the different source processes and of independence of the inter-arrivals and the sizes of each source process, no longer serve as good approximations in describing the behavior of the source processes. (A small amount of queueing, which almost always occurs, does not invalidate the approximation.)

Consider two packets, $j = 19$ and $j = 20$. Suppose packet 20 waits in the queue for packet 19 to be transmitted. The two are back-to-back on the link, which means, because the arrival time is the first moment of transmission, that $t_{19}$ is the time to put the bits of packet 19 on the link, which is equal to $q_{19}/\ell$, where $\ell$ is the link speed. For example, at $\ell = 100$ mbps, the time for a 1500 byte (12000 bit) packet is 120 $\mu$. So given $q_{19}$ we know $t_{19}$ exactly. Queueing can occur on routers further upstream than the link-input router and affect the assumptions as well.

The arrival times of the packets on the link, $a_j$, are the departure times of the packets from the queue. The departure times are the arrival times at the queue plus the time spent in the queue. If there are no other packets in the queue when packet $j$ arrives, then $a_j$ is also the arrival time at the queue. Suppose queueing is first-in-first-out. Then the order of the arriving packets at the queue is the same as the order of departing packets from the queue, so $q_j$ is also the packet size process for the arrivals at the queue.

The effect of queueing on $q_j$ is simple. Because the queueing does not alter the $q_j$, the statistical properties of the $q_j$ are unaffected by the queueing; in particular, their limit of independence is not altered..

But statistical theory for the departure times from a queue is not developed well enough to provide much guidance for the affect of queueing on the statistical properties of $t_j$ and $p_i$. However, the properties of the extreme case are clear. If $m$ is so large that the queue never drains, then the $t_j$ are equal to $q_j/\ell$, so the $t_j$ take on the statistical properties of $q_j$. Since the $q_j$ tend to independence, the $t_j$ eventually go to independence, so there is no long-range dependence. A Poisson process is a renewal process, a point process with independent inter-arrivals, with the added property that the marginal distribution of the inter-arrivals is exponential. The extreme $t_j$ process is a renewal process but with a marginal distribution proportional to that of the packet sizes. The extreme $p_i$ is the count process corresponding to the $t_j$ renewal process; this implies the coefficient variation of $p_i$ is a constant, so the decrease like $1/\sqrt{m}$ prescribed by the superposition theory is arrested, and it implies the $p_i$ are independent, so there is no long-range dependence. We do not expect to see the extreme case in our empirical study, but it does provide at least a point of attraction.

## VI. EMPIRICAL STUDY: INTRODUCTION

The superposition theory and the heuristic discussion of the effect of upstream queueing provide hypotheses about the statistical properties of the inter-arrivals $t_j$, the sizes $q_j$, and the counts $p_i$. We carried out extensive empirical studies to investigate the validity of the hypotheses [10], [19], [25].

In the early 1990s, Internet researchers put together a comprehensive measurement framework for studying the characteristics of packet traffic that allows not just statistical study of traffic, but performance studies of Internet engineering designs, protocols, and algorithms [26], [27]. The framework consists of capturing the headers of all packets arriving on a link and time-stamping the packet, that is, measuring the arrival time, $a_j$. The result of measuring over an interval of time is a *packet trace*. Packet trace collection today enjoys a very high degree of accuracy and effectiveness for traffic study [28], [29].

| Trace Group | Number | Link | $\log(c)$ |
|---|---|---|---|
| AIX1(90sec) | 23 | 622mbps | 13.09 |
| AIX2(90sec) | 23 | 622mbps | 13.06 |
| COS1(90sec) | 90 | 156mbps | 10.83 |
| COS2(90sec) | 90 | 156mbps | 10.81 |
| NZIX(5min) | 100 | 100 mbps | 10.75 |
| NZIX7(5min) | 100 | 100 mbps | 9.60 |
| NZIX5(5min) | 100 | 100 mbps | 8.66 |
| NZIX6(5min) | 100 | 100 mbps | 7.85 |
| NZIX2(5min) | 100 | 100 mbps | 7.32 |
| NZIX4(5min) | 100 | 100 mbps | 7.17 |
| BELL(5min) | 500 | 100 mbps | 6.97 |
| NZIX3(5min) | 100 | 100 mbps | 6.54 |
| BELL-IN(5min) | 500 | 100 mbps | 5.98 |
| BELL-OUT(5min) | 500 | 100 mbps | 5.94 |
| NZIX1(5min) | 100 | 100 mbps | 4.42 |

TABLE I

LINK: NAME INCLUDING LENGTH OF TRACES • NUMBER:
NUMBER OF TRACES • LINK: SPEED • $\log(c)$: LOG BASE 2
AVERAGE NUMBER OF ACTIVE CONNECTIONS

We put together a very large database of packet traces measuring many Internet links whose speeds range from 10 mbps to 2.5 gbps, and we built S-Net, a software system, based on the S language for graphics and data analysis, for analyzing very large packet header databases [20]. We put the database and S-Net work to study the multiplexing hypotheses.

For each studied trace, which covers a specific block of time on a link, we compute $a_j$, $t_j$, $q_j$, and 100-ms $p_i$. We also need a summary measure of the magnitude of multiplexing for the trace. At each point in time over the trace, the measure is the number of active connections. The summary measure, $c$, for the whole trace is the average number of active connections over all times in the trace. Here, we describe some of the results of one of our empirical investigations in which we analyzed 2526 header packet traces, 5 min or 90 sec in duration, from 6 Internet monitors measuring 15 links ranging from 100 mbps to 622 mbps [19]. Table I shows information about the traces. Each row describes the traces for one link. The first column gives the trace group name: the trace length is a part of each name. Column 2 gives the number of traces. Column 3 gives the link speed. Column 4 gives the mean of the log base 2 of $c$ for the traces of the link.

Consider each packet in a trace. Arriving after it is a back-to-back run of $k$ packets, for $k = 0, 1, \ldots$; each packet in the run is back-to-back with its predecessor. If packet 19 has a back-to-back run of 3 packets, then packet 20 is back-to-back with 19, 21 is back-to-back with 20, 22 is back-to-back with 21, but 23 is not back-to-back with 22. The percent of pack-

ets with back-to-back runs of $k$ or more is a measure of the amount of queueing on the link-input router. We studied this measure for many values of $k$. We need such study to indicate when the network is likely pushing back on the attraction to Poisson and independence.

Figure 1 graphs the percent of packets whose back-to-back runs are 3 or greater against $\log(c)$. Each point on the plot is one trace. Each of the 15 panels contains the points for one link. The panels are ordered, left to right and bottom to top, by the means of the $\log(c)$ for the 15 links, given in column 4 of Table I.

Figure 1, and others like it for different values of $k$, show that only four links experience more than minor queueing — COS1, COS2, AIX1, and AIX2 — so we would not expect to see significant push-back except at these four. However, queueing further upstream than the link-input router can affect the traffic properties as well, but without creating back-to-back packets, so we reserve final judgment until we see the coming analyses.

Figure 1 also provides information about the values of $c$. Since the mean of $\log(c)$ increases left to right and bottom to top, the distribution shifts generally toward higher values in this order. The smallest $c$, which appears in the lower left panel, is 5.9 connections; the largest, which appears in the upper right panel, is 16164 connections.

## VII. EMPIRICAL STUDY: FSD AND FSD-MA(1) MODELS

In this section we introduce two very simple classes of stationary time series models [25], one a subclass of the other, that we found provide excellent fits to the inter-arrivals $t_j$, the sizes $q_j$, and the counts $p_i$ for the 2526 traces. The models are parametric. One of the parameters determines the amount of dependence. At low values of the parameter, the series has substantial autocorrelation and is long-range dependent. As the parameter increases, the amount of dependence decreases. At the largest value of the parameter, the series is independent. Other parameters determine the marginal distribution of the series and therefore the coefficient of variation. By fitting the models to each trace, we can study the multiplexing gains by studying the changing values of the parameters across the traces, and relating the changes to the average active connection load $c$ of the traces.

The two model classes are fractional sum-difference (FSD) models and FSD-MA(1) models [25]. FSD models have two additive components:
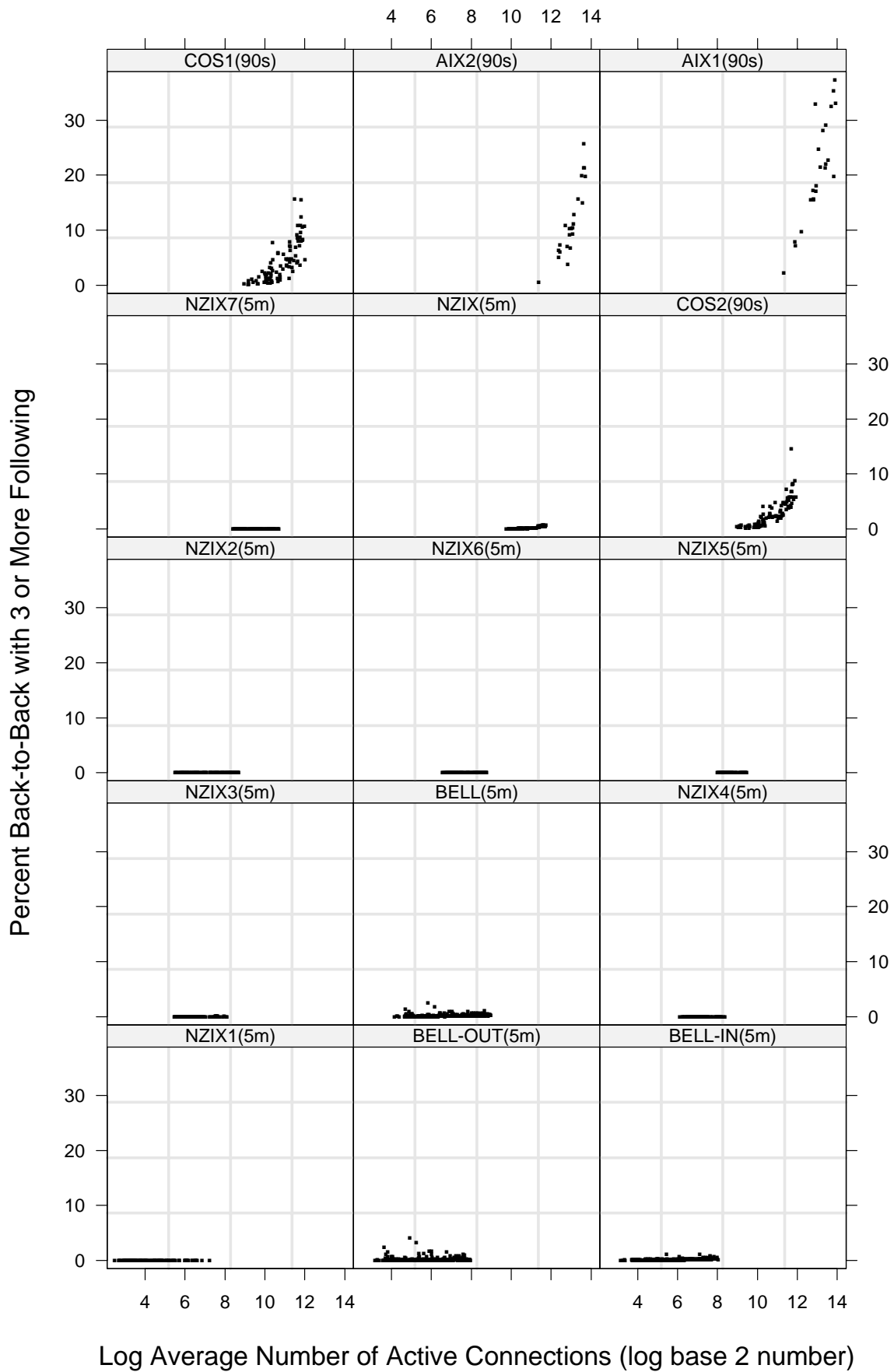
Fig. 1. The percent of packets with back-to-back runs of 3 or greater is plotted against $\log(c)$.

a simple fractional ARIMA and white noise. MA(1) refers to a first-order moving average [30]. FSD-MA(1) models replace the white noise of the FSD model with an MA(1). Since white noise is a special case of an MA(1), the FSD models are a subclass of the FSD-MA(1) models. As we will see, the names "transformation-Gaussian FSD models" and "transformation-Gaussian FSD-MA(1) models" would convey more information about the nature of the models, but for simplicity we will use the shorter names.

Suppose $x_u$ for $u = 1, 2, \ldots$ is a stationary time series with a marginal cumulative distribution function $F(x; \phi)$ where $\phi$ is a vector of unknown parameters. For example, $F(x; \phi)$ might be log-normal or Weibull. Let $z_u = T(x_u; \phi)$ be a transformation of $x_u$ such that the marginal distribution of $z_u$ is normal with mean 0 and variance 1. If $G^{-1}(r)$ is the quantile with probability $r$ of $z_u$, then $T(x_u; \phi) = G^{-1}(F(x_u; \phi))$. If $x_u$ is log-normal and the vector $\phi$ consists of the mean $\mu$ and variance $\sigma^2$ on the log scale, then $T(x_u; \phi) = (\log(x_u) - \mu)/\sigma$.

Next we suppose that $z_u$ is a Gaussian time series, that is, the joint distributions of all finite subsets of the time series are multivariate normal.

Let

$$z_u = \sqrt{1 - \theta}\, s_u + \sqrt{\theta}\, n_u,$$

where $s_u$ and $n_u$ are independent of one another and each has mean 0 and variance 1. $n_u$ is white noise, that is, an uncorrelated time series. $s_u$ is a fractional ARIMA (FARIMA) model [31]

$$(I - B)^d s_u = \epsilon_u + \epsilon_{u-1}$$

where $Bs_u = s_{u-1}$, $0 < d < 0.5$, and $\epsilon_u$ is white noise with mean 0 and variance

$$\sigma_\epsilon^2 = \frac{(1 - d)\Gamma^2(1 - d)}{2\Gamma(1 - 2d)}$$

to make the variance of $s_u$ equal to 1.

$z_u$ is an FSD model. We coined this term because the model for $z_u$ can be written as a combination of fractional and summation difference operators acting on $z_u$ and on two white noise series:

$$(I - B)^d z_u = (I + B)\epsilon_u + (I - B)^d n_u.$$

These models are to FARIMA models what the very simple and widely applicable IMA(1,1) models are to ARIMA models [30]; the IMA(1,1) models can be written as

$$(I - B)z_u = (I + B)\epsilon_u + (I - B)n_u.$$

Generalizations of this latter model have been named *sum-difference models* [32].

The FSD-MA(1) model is

$$z_u = \sqrt{1 - \theta}\, s_u + \sqrt{\theta}\, n_u,$$

similar to the FSD model, but where $n_u$ instead of white noise is a first order moving-average

$$n_u = \zeta_u + \beta\zeta_{u-1},$$

where $\zeta_u$ is Gaussian white noise with mean 0 and variance $(1 + \beta^2)^{-1}$, which makes the variance of $n_u$ equal to 1. If $\beta = 0$, the moving-average component is white noise so the model is simply an FSD. We need the above restriction $d > 0$. If $d = 0$, the model is not identifiable because $z_u$, whose model has two parameters, is a first order moving average with variance 1, which has one parameter.

Suppose $z_u$ is an FSD-MA(1) model. Let $r_z(k), r_s(k)$ and $r_n(k)$ be the autocorrelation functions of of $z_u$, $s_u$, and $n_u$, respectively, for lags $k = 0, 1, 2, \ldots$ . Because $d > 0$, $s_u$ is long-range dependent, and $r_s(k)$ falls off like $k^{2d-1}$ and increases at all positive lags as $d$ increases. $r_n(k) = \beta(1 + \beta^2)^{-1}\{k = 1\}$ where $\{k = 1\}$ is 1, if $k = 1$, and is 0 if $k > 1$. Thus

$$r_z(k) = (1 - \theta)r_s(k) + \theta\beta(1 + \beta^2)^{-1}\{k = 1\}.$$

As $\theta \to 1$, the long-range dependent component $\sqrt{1 - \theta}\, s_u$ contributes less and less variation to $z_u$. Finally, when $\theta = 1$, $z_u$ is white noise if $\beta = 0$, and is a first-order moving average otherwise.

The power spectrum of $z_u$ is

$$p_z(f) = (1-\theta)\sigma_\epsilon^2\frac{4\cos^2(\pi f)}{\left(4\sin^2(\pi f)\right)^d} + \theta\frac{1 + \beta^2 + 2\beta\cos(2\pi f)}{1 + \beta^2}$$

for $0 \leq f \leq 0.5$. The frequency $f$ has units cycles/inter-arrival for $t_j$, cycles/packet for $q_j$, and cycles/interval-length for $p_i$. $p_z(f)$ decreases monotonically as $f$ increases. Because $d > 0$, the term $\sin^{-2d}(\pi f)$ goes to infinity at $f = 0$, so if $\theta < 1$, no matter how close $\theta$ gets to 1, $p_z(f)$ gets arbitrarily large near $f = 0$, but its ascent begins closer and closer to 0 as $\theta$ gets closer to 1.

Figure 2 shows the power spectra for 16 FSD-MA(1) models. For each panel, the spectrum is evaluated at 100 frequencies, equally spaced on a log base 2 scale from $-13$ to $-1$. The value of $d$ in all 16 cases is $0.41$. $\theta$ varies from 0.39 to 0.99 by 0.2 as we go left
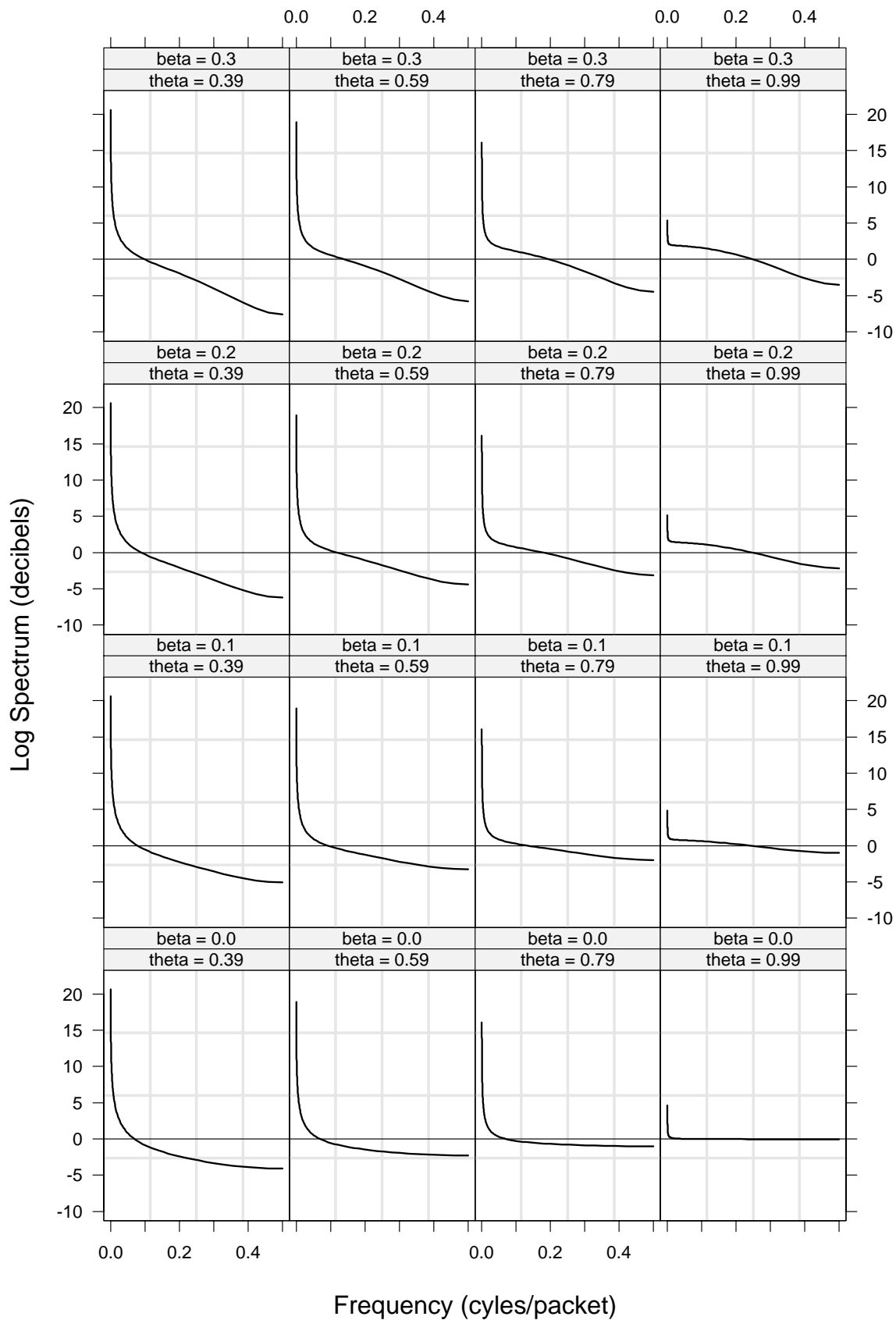
Fig. 2. The log power spectrum of an FSD-MA(1) time series is plotted against frequency for different values of $\theta$ and $\beta$.

to right through the columns. $\beta$ varies from 0 to 0.3 by 0.1 as we go from bottom to top through the rows.

So the bottom row shows spectra for the FSD model, while the other rows show the spectra for FSD-MA(1) models with positive $\beta$.

For all panels, there is a rapid rise as $f$ tends to 0, and an overall monotone decrease in power as the frequency increases from 0 to 0.5. This is a result of the persistent long-range dependence. But for each row, as $\theta$ increases, the fraction of low-frequency power decreases, and the fraction of high-frequency power increases. In the bottom row, as $\theta$ increases, the spectrum at frequencies away from 0 shows a distinct flattening, tending toward the flat spectrum of white noise. In the remaining rows, the spectra, away from 0, tend toward that of a gently sloping curve, the spectrum of an MA(1).

We found that the 100-ms packet counts, $p_i$, and the packet sizes, $q_j$, for all but a few of the 2526 traces, are very well fitted by an FSD model. $t_j$ is also typically well fitted by either an FSD model or an FSD-MA(1); for the traces of some links, an FSD-MA(1) model with a positive $\beta$ is clearly required as $c$ gets large, the result of the network pushing back on the attraction to Poisson and independence by upstream queueing.

The estimation of the parameters, especially of $d$, needs considerable care. But an essential part of the study was visualization tools that validated the resulting fitted models. The estimation and modeling checking is discussed in detail elsewhere [25].

## VIII. EMPIRICAL STUDY: PACKET COUNTS

The superposition theory predicts that the coefficient of variation of the $p_i$ should decrease like $1/\sqrt{c}$. Figure 3 graphs the log of the coefficient against $\log(c)$. The theory predicts a slope of $-0.5$; the least squares line with slope $-0.5$ is shown on each panel. The rate of decline of the log coefficients is certainly consistent with a value of $-0.5$. At some sites, the decline is somewhat faster and at others, it is slower. Interestingly, the decline has not been altered by back-to-back occurrence, as predicted by the heuristics for the effect of upstream queueing, even for AIX1 and AIX2 which have the largest back-to-back percents. This presumably happened in part because the aggregation interval length is 100 ms; had we used a smaller interval, an effect might have been detected.

The $p_i$ do not have a normal marginal until $c$ gets large. A log-normal marginal does much better. Let $p_i^*$ be $\log(1 + p_i)$ normalized to have mean 0 and variance 1. An assumption of a Gaussian process for $p_i^*$ is a reasonable approximation for much smaller $c$. We found that an FSD model fitted the $p_i^*$ extremely well,

except for a small fraction of intervals with low $c$ where oscillatory effects of Internet transport protocols broke through and created spikes in the power spectrum. Even in these cases, the model serves as an excellent summary of the amount of long-range dependence in the correlation structure.

Estimates of $d$ vary by a small amount across the traces and showed no dependence on $c$. The medians for the 15 links vary from 0.39 to 0.45 and their mean is 0.41. Estimates of $\theta$ also show no dependence on $c$; the mean of the 15 medians of $\theta$ is 0.53. Thus the $p_i^*$ spectra look like the spectrum in column 2 and row 1 in Figure 2. The stability of the correlation structure of $p_i^*$ is consistent with the superposition theory, which stipulates that the correlation structure of the $p_i$ does not change with $c$. The heuristics for the effect of upstream queueing suggest that the autocorrelation should be changed by a large amount of upstream queueing, but the effect does not appear to occur even for AIX1 or AIX2, where the occurrence of back-to-back packets is the greatest. As with the coefficient of variation, it is possible an effect would be seen for interval lengths less than 100 ms.

## IX. EMPIRICAL STUDY: PACKET SIZES

A reasonable summary of the marginal distribution of the $q_j$ is an atom at the minimum size of 40 bytes, an atom at the maximum size of 1500 bytes, an atom at 576 bytes, and continuous uniform from 40 bytes to 1500 bytes. Quantile plots [33] showed that the marginal distribution did not change appreciably with $c$, as predicted by the superposition theory, but did change appreciably across the 15 links. For example, if a link has traffic in a single direction from hosts with a preponderance of clients downloading web pages, then the frequency of 40 byte packets is greater and the frequency of 1500 byte packets less than for a link where the preponderance of hosts are serving web pages.

We do not transform the $q_j$ to a normal marginal for our FSD modeling because the transformation would not be invertible. For analysis purposes, we treat the $q_j$ as is, without transformation; this amounts to a second moment analysis, but it will provide adequate insight because the correlation coefficient is still a reasonable summary of dependence for such discrete-continuous data.

We found that an FSD model provided an excellent fit to the $q_j$. A combination of theory and empirical study show that $d$ remains constant with $c$, and the estimate of the single value came out to 0.42, very close to the 0.41 for the $p_i^*$. For simplicity, we could
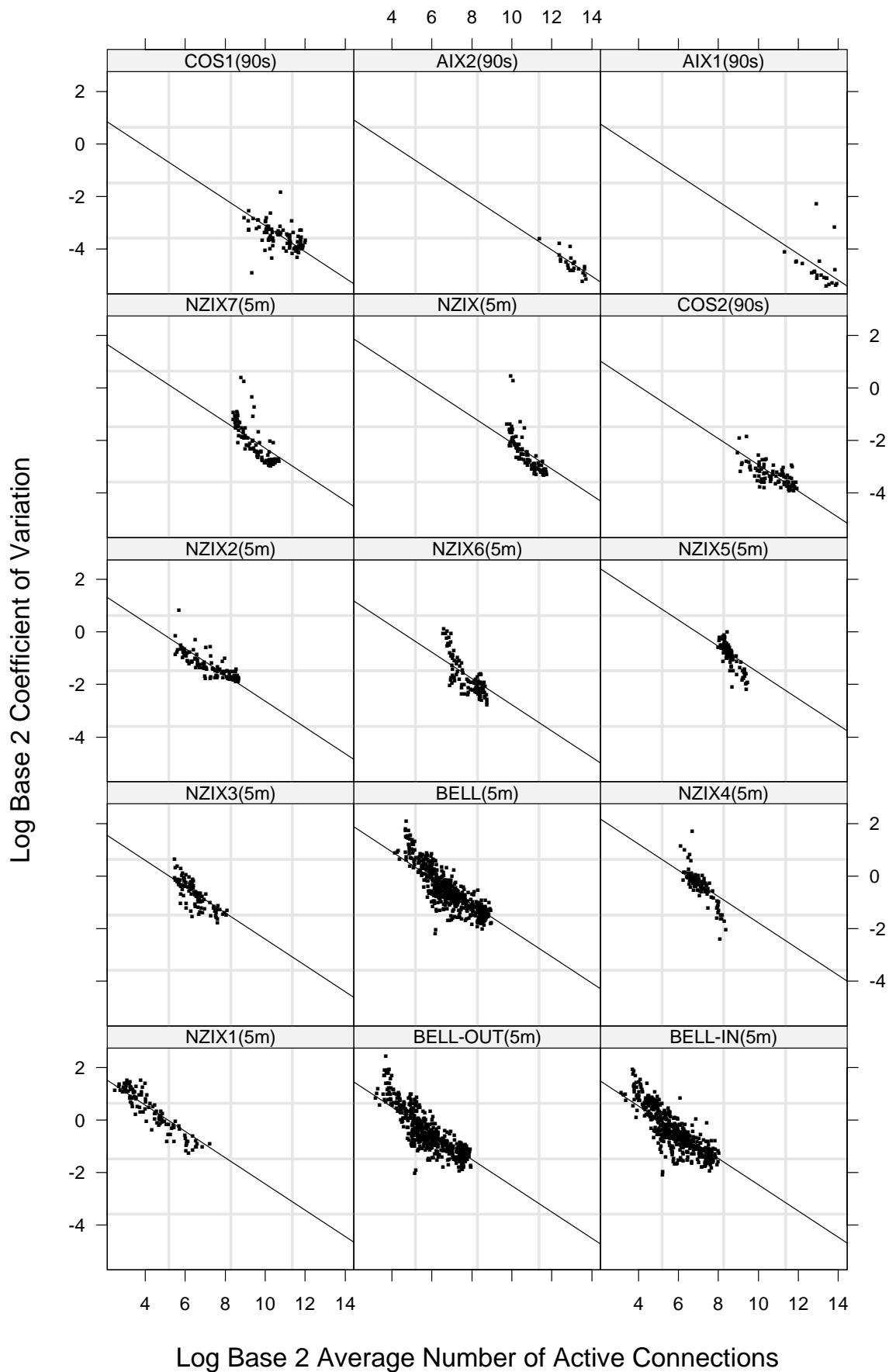
Fig. 3. Log coefficient of variation of the 100-ms packet counts is plotted against $\log(c)$.
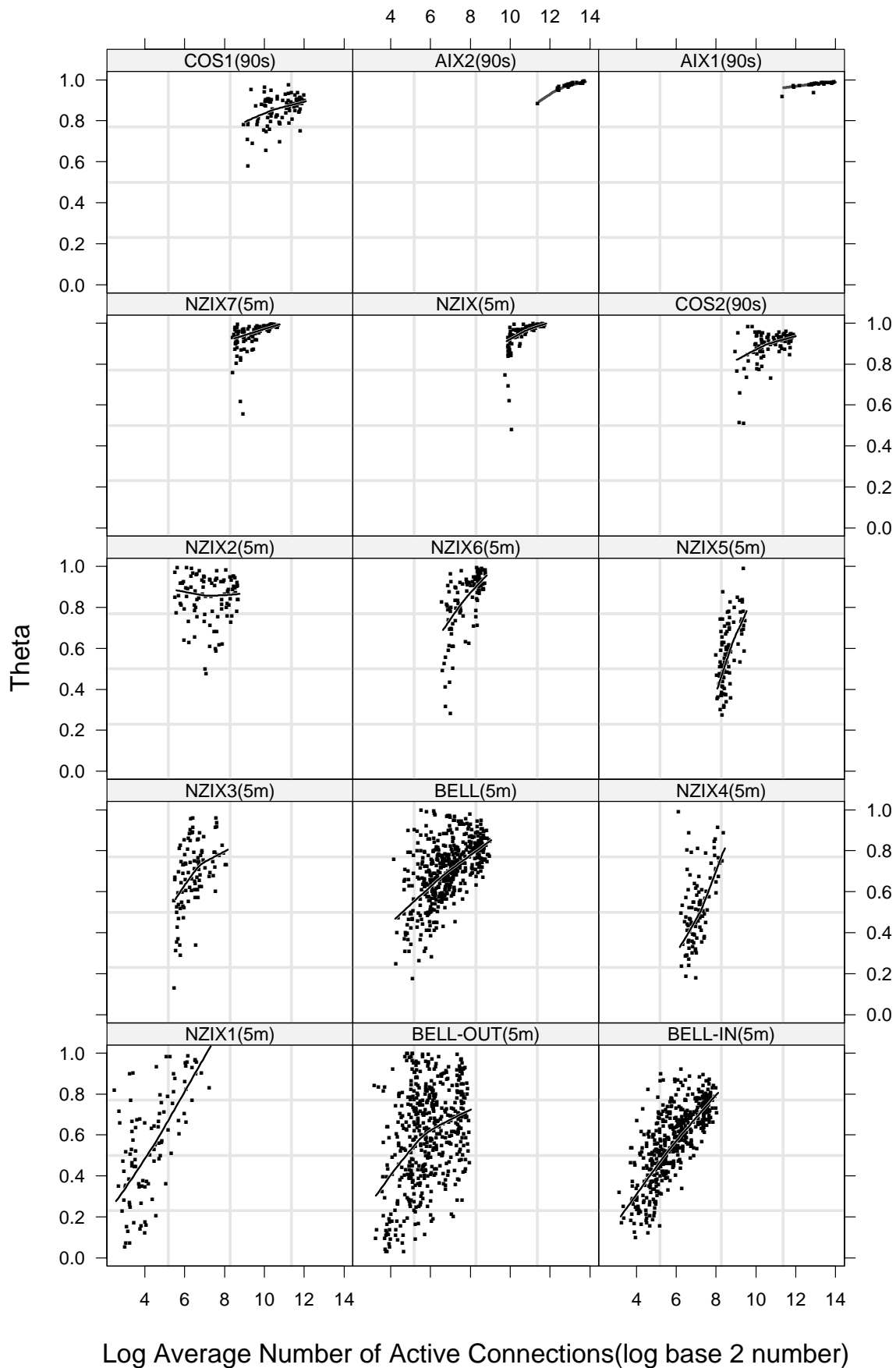
Fig. 4. An estimate of $\theta$ for the packet sizes is plotted against $\log(c)$.

not resist using a value of 0.41 for the $q_j$, the same as the estimate of $d$ for the $p_i$. We fixed $d$ to this value and estimated $\theta$.

Figure 4 plots the estimates of $\theta$ against $\log(c)$. The smooth curve on each panel is a loess fit using robust locally linear fitting and a smoothing parameter of 1 [33]. Loess is a nonparametric procedure that puts curves through data by a moving local polynomial fitting procedure, the same in spirit, but not in detail, as a moving average smoothing a time series. The overall result in Figure 4 is that $\theta$ goes to 1 with $c$, so the spectrum changes as shown in the bottom row of Figure 2. Thus the $q_j$ tend toward independence as prescribed by the the superposition theory. An increase in the percent of back-to-back packets with $c$ for the COS1, COS2, AIX1, and AIX2 links does not alter the increase in $\theta$, which is consistent with the heuristics for the effect of upstream queueing.

## X. EMPIRICAL STUDY: INTER-ARRIVALS

We found, using Weibull quantile plots, that the marginal distribution of the inter-arrivals is well approximated by the Weibull distribution across all values of $c$. The back-to-back packets result in deviations from the Weibull, but because packet sizes vary by a factor of $1500/40 = 37.5$, the deviations are spread across the distribution, and overall the approximation remains excellent, even for the traces with the largest occurrence of back-to-back packets. The Weibull has two parameters: $\alpha$, a scale parameter, and $\lambda$ a shape. When $\lambda$ is 1, the Weibull is an exponential, the inter-arrival distribution of a Poisson process. When $\lambda < 1$, the tail is heavier than that of the exponential.

Figure 5 plots estimates of the Weibull shape parameter, $\lambda$, against $\log(c)$. The smooth curve on each panel is a loess fit with robust locally linear fitting and a smoothing parameter of 1. The overall result is that the shape estimates are less than 1, and as $c$ increases, the shape tends toward 1. Consider the 5 links with the largest mean $\log(c)$ — NZIX, COS2, COS1, AIX2, and AIX1. Almost all of the values of $c$ exceed $2^{10}$, but few values for the remaining sites do so. For these top five, most estimates of $\lambda$ exceed 0.9. For the remaining, most estimates are below 0.8. The top five appear to have a limit slightly less than 1; the back-to-back packets exert just enough influence to keep the estimates slightly below 1, but this is a small matter since a Weibull with shape of 0.95 is exceedingly close to exponential.

Because the $t_j$ have a Weibull marginal, the transformation that takes them to normality is $T(t_j; \phi) =$

$G^{-1}(F(t_j; \phi))$ where $F$ is the Weibull cumulative distribution function and $\phi$ is the vector of parameters $\alpha$ and $\lambda$. Because $\lambda$ changes, the transformation changes, but the change is not large and we found the transformations are well approximated by a single transformation, the sixth root of $t_j$. So for simplicity we used $t_j^* = t_j^{1/6}$.

We found that an FSD model or an FSD-MA(1) model provided an excellent fit to $t_j^*$ except for a small fraction of intervals with low $c$ where oscillatory effects of Internet transport protocols broke through and created spikes in the power spectrum. Even in these cases, the model serves as an excellent summary of the amount of long-range dependence in the correlation structure.

Theoretical results show that the value of $d$ for the $t_j$, or for monotone transformations of $t_j$ such as $t_j^*$, is the same as that for the $p_i^*$, so the estimate of $d$ for the $t_j^*$ was taken to be 0.41, that for $p_i^*$. This was done rather than estimating $d$ from the $t_j^*$ because, when $\theta$ gets close to 1, the long-range dependent component accounts for such a small fraction of the variation in the $t_j^*$ that $d$ is poorly estimated.

We fitted an FSD-MA(1) with $d = 0.41$ to the 2526 traces. Figure 6 graphs estimates of $\beta$ against $\log(c)$. The smooth curve on each panel is a loess fit with robust locally linear fitting and a smoothing parameter of 1. 1.2% of the estimates are less than $-0.4$ and are not shown on the plot. Our model checking showed that the MA(1) component was important for producing a good fit for the largest values of $c$ at NZIX7, NZIX, AIX1, and AIX2. The latter two sites show a large back-to-back occurrence, but not the first two. However, queueing upstream from the link-input router can affect the inter-arrivals without introducing back-to-back packets. In other words, our measure of back-to-back packets in Figure 1 does not tell the whole story of upstream queueing.

For $\beta \leq 0.1$, $n_u = \zeta_u + \zeta_{u-1}$ is nearly white noise. When $\beta = 0.1$, the variance of $\zeta_u$ is $1/(1 + .1^2) = 0.990$, so $n_u$, whose variance is 1, is very close to white noise. But when $\beta = 0.3$, the variance of $\zeta_u$ is 0.917, so $n_u$ contains significant correlated variation. It is only at NZIX7, NZIX, AIX1, and AIX2 that $\beta$ is reliably above 0.1, getting as high as 0.3. At the other links, $\beta$ is small enough, taking the greater variability of estimates as $c$ decreases into account, that it is reasonable to omit the MA(1) component, that is, using just an FSD model. In particular, at COS1 and COS2, $\beta$ is small.

Figure 7 graphs $\theta$ against $\log(c)$. The smooth curve on each panel is a loess fit with robust locally linear
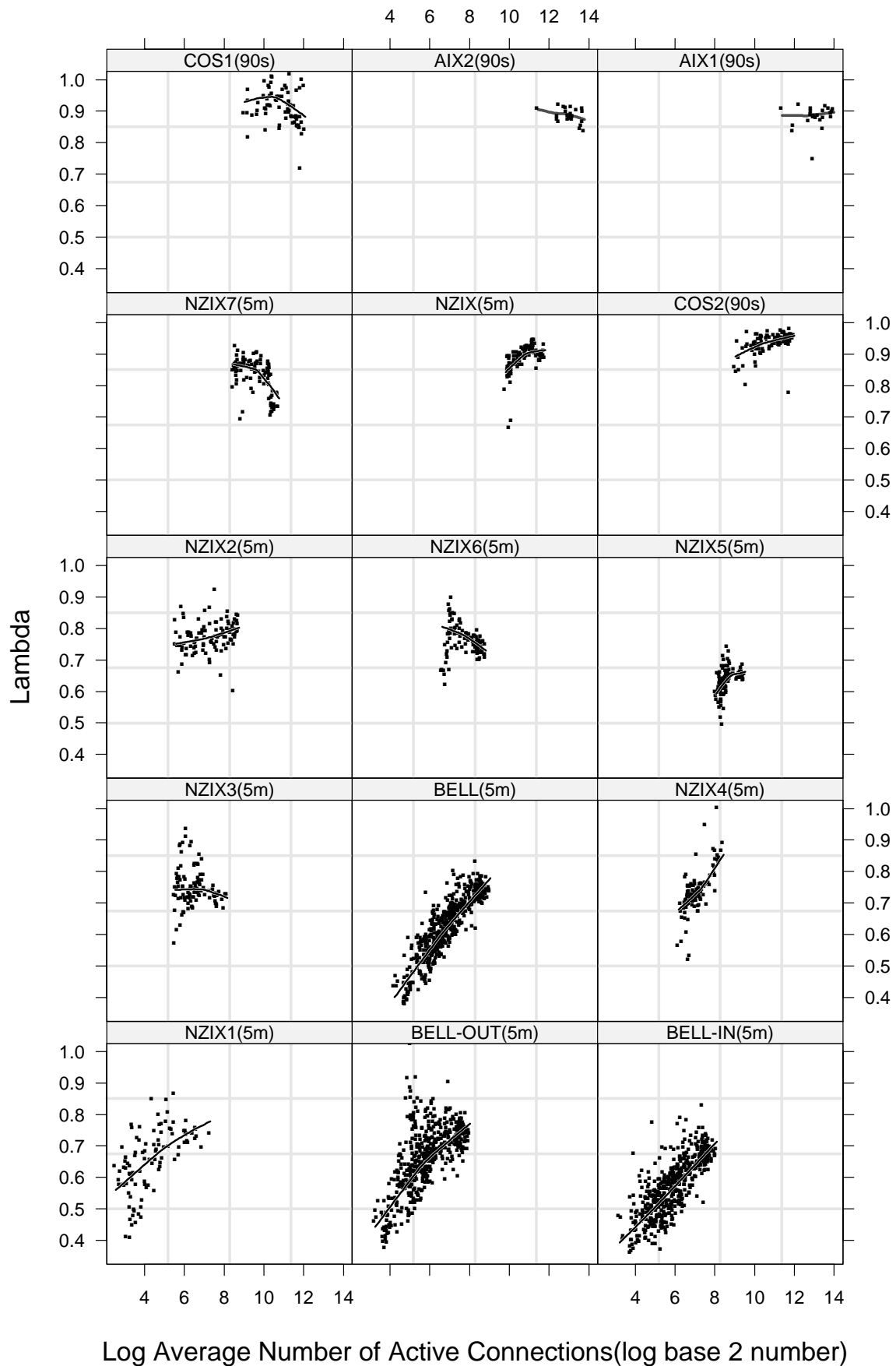
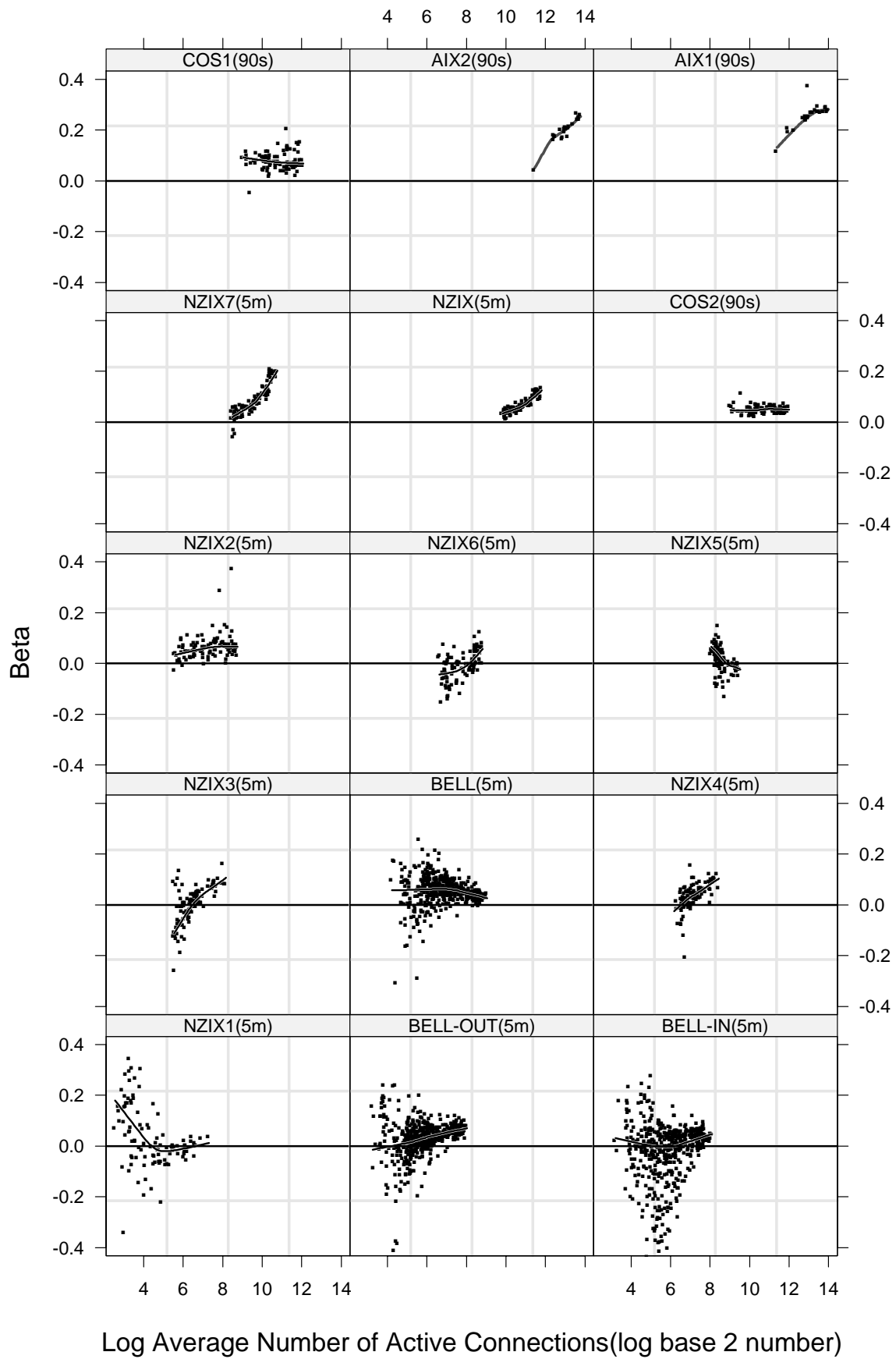Fig. 5. An estimate of $\lambda$ for the inter-arrivals is plotted against $\log(c)$.

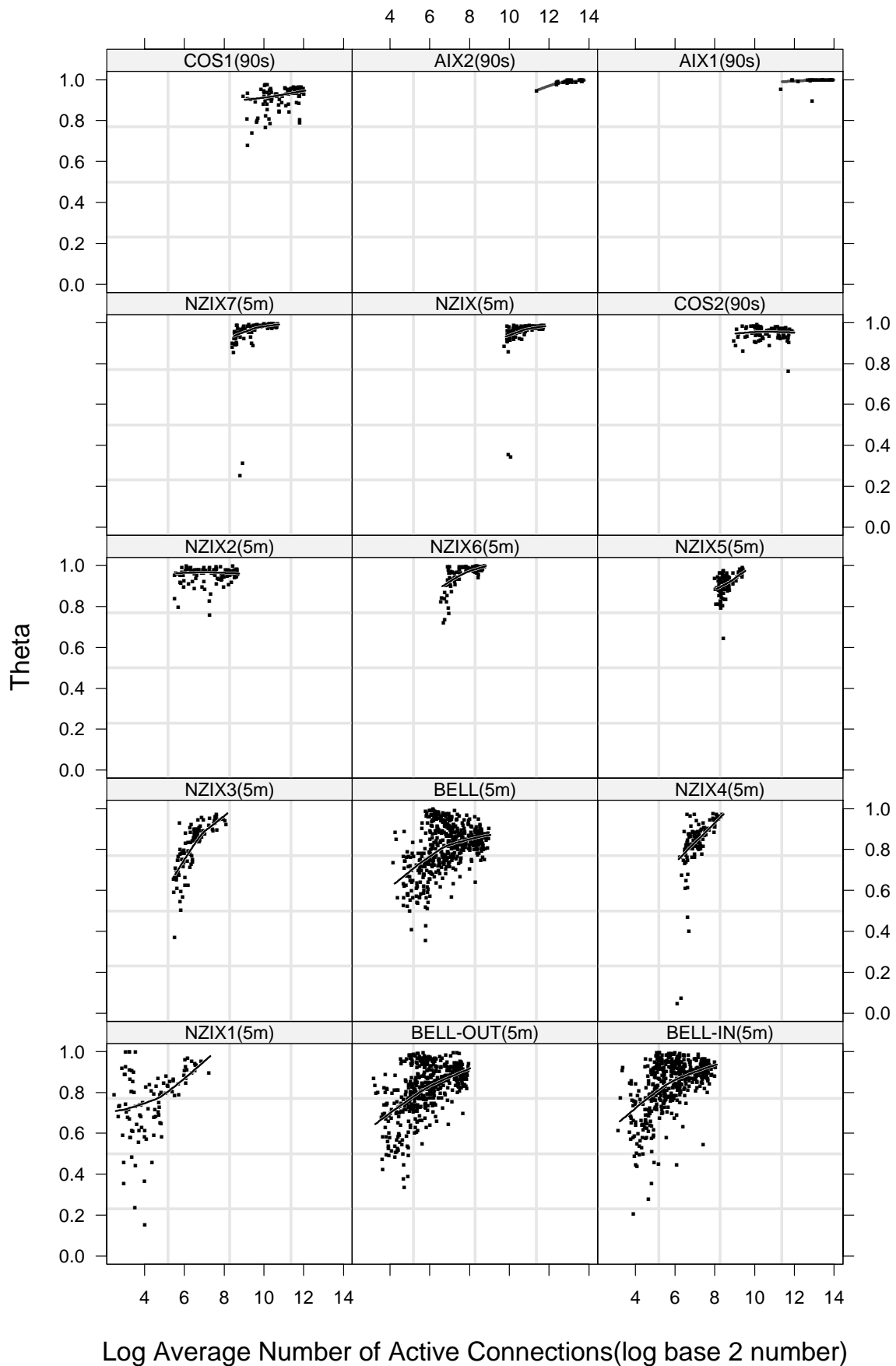Fig. 6. An estimate of $\beta$ for the inter-arrivals is plotted against $\log(c)$.

Fig. 7. An estimate of $\theta$ for the inter-arrivals is plotted against $\log(c)$.

fitting and a smoothing parameter of 1. The overall result is that $\theta$ goes to 1 with $c$. The long-range dependence of $t_j^*$ dissipates, tending either to short-range dependence, an MA(1), or to independence. Thus all panels of Figure 2 convey the behaviors of the power spectra of the $t_j^*$.

These results for $\lambda$, $\beta$, and $\theta$ are consistent with the superposition theory and the heuristics for the effect of upstream queueing. Multiplexing creates an attraction to Poisson in the $t_j$; $\lambda$ and $\theta$ tend toward 1 as the theory prescribes. But the network succeeds in pushing back in some cases, keeping $\lambda$ slightly less than 1, and causing values of $\beta$ for some links that indicate short-term dependence.

## XI. OPEN LOOP GENERATION OF PACKET TRAFFIC

The FSD models fitted to the sizes and inter-arrivals can be used for open-loop generation of synthetic traffic for simulation studies. The inter-arrival marginal is Weibull; the parameters are $\alpha$ and $\lambda$. The packet size marginal has atoms at specific packet sizes and has a continuous part that is uniform between 40 bytes and 1500 bytes; the parameters are the probabilities at the atoms. The inter-arrivals are generated by Gaussian FSD variables with d = 0.41 transformed to the Weibull marginal; the parameter is $\theta_t$. The packet sizes are generated by Gaussian FSD variables with d = 0.41 transformed to the discrete-continuous marginal; the parameter is $\theta_q$. $\alpha$, $\lambda$, $\theta_t$, and $\theta_q$ change with $c$ according to certain models to reflect the multiplexing gains, so only $c$ is specified to carry out generation.

## XII. THERE ARE MULTIPLEXING GAINS

The results here show that an increasing number of simultaneous active connections causes a dramatic change in the statistical properties of packet traffic on an Internet link. Starting at low connection loads on an uncongested link, packet arrivals tend toward Poisson and packet sizes tend toward independence as the load increases. A component of long-range dependence is retained in each of these variables, but the effect of the component gets increasingly small. Packet counts have a stable autocorrelation structure that does not change with the load, but the standard deviation of the counts relative to the mean gets small, so the counts become smooth. The network pushes back on this attraction to Poisson and independence through upstream queueing, which also increases with the connection load; very short term autocorrelation can develop in the inter-arrivals,

and their marginal changes toward the distribution of packet sizes divided by the link speed. On a link with a sufficiently large speed that the increasing connection load can bring the traffic to Poisson and independence before substantial upstream queueing occurs, the onset of queueing does not resurrect the long-range dependence. All this means that the burstiness of traffic, once thought to pervade the whole Internet, dissipates with the connection load. There are multiplexing gains.

Inspired by these results on multiplexing gains, theoretical and empirical studies have now demonstrated that queueing on an Internet device tends to that of Poisson arrivals and independent sizes as the load increases, just as one would expect [10], [21]. This means that if a link speed is sufficiently large, queueing distributions relative to the bit/rate of the traffic get dramatically smaller.

The foundations of traffic analysis and modeling should reflect these results. The dramatic change in the statistical properties with the connection load makes clear that the load needs to play a central role in analysis and modeling. Theory must reflect the load. Empirical study must encompass a range of packet traces from small loads to large.

The results have important implications for Internet device engineering and Internet traffic engineering. On links with low speeds, at the edges of the Internet close to the user hosts, connection loads cannot get large, and traffic remains highly bursty. But on links with high speeds, toward the core of the Internet and carrying traffic made up of large numbers of connections, the traffic can be close to Poisson and independence, so the burstiness is gone. Engineering studies that are meant to apply to the Internet as a whole, and that use synthetic or live packet traffic to assess performance, need to consider packet traces varying across a wide range of link speeds and connection loads. Many issues of Internet engineering need to be revisited to determine how protocols, algorithms, device design, network design, and network provisioning should change to reflect the effect of the changing statistical properties of the traffic with the connection load.

## XIII. ACKNOWLEDGEMENTS

## References

[1] W. Leland, M. Taqqu, W. Willinger, and D. Wilson, "On the Self-Similar Nature of Ethernet Traffic," *IEEE/ACM Transactions on Networking*, vol. 2, pp. 1–15, 1994.

[2] Vern Paxson and Sally Floyd, "Wide-Area Traffic: The Failure of Poisson Modeling," *IEEE/ACM Transactions on Networking*, vol. 3, pp. 226–244, 1995.

[3] M. E. Crovella and A. Bestavros, "Self-Similarity in World Wide Web Traffic: Evidence and Possible Causes," *ACM SIGMETRICS*, pp. 160–169, 1996.

[4] A. Feldman, A. A. Gilbert, and W. Willinger, "Data Networks as Cascades: Explaining the Mulifractal Nature of Internet WAN Traffic," in *Proceedings ACM SIGCOMM*, 1998, pp. 42–55.

[5] Rudolf H. Riedi, Matthew S. Crouse, Vinay J. Ribeiro, and Richard G. Baraniuk, "A Multifractal Wavelet Model with Application to Network Traffic," *IEEE Transactions on Information Theory*, vol. 45, no. 3, pp. 992–1019, 1999.

[6] B. B. Mandelbrot, "Long-Run Linearity, Locally Gaussian Processes, H-Spectra and Infinite Variances," *International Economic Review*, vol. 10, pp. 82–113, 1969.

[7] Sally Floyd and Vern Paxson, "Why We Don't Know How to Simulate the Internet," Tech. Rep., LBL Network Research Group, 1999.

[8] K. Park, G. Kim, and M. Crovella, "On the Relationship Between File Sizes, Transport Protocols, and Self-Similar Network Traffic," in *Proceedings of the IEEE International Conference on Network Protocols*, 1996.

[9] 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*, vol. 5, pp. 71–86, 1997.

[10] J. Cao, W. S. Cleveland, D. Lin, and D. X. Sun, "On the Nonstationarity of Internet Traffic," *ACM SIGMETRICS*, pp. 102–112, 2001.

[11] J.B. Gao and I. Rubin, "Multiplicative Multifractal Modeling of Long-Range-Dependent Network Traffic.," *International Journal of Communications Systems*, vol. 14, pp. 783–201, 2001.

[12] Ashok Erramilli, Onuttom Narayan, and Walter Willinger, "Experimental Queueing Analysis with Long-Range Dependent Packet Traffic," *IEEE/ACM Transactions on Networking*, vol. 4, pp. 209–223, 1996.

[13] D. D. Botvich and N. G. Duffield, "Large Deviations, the Shape of the Loss Curve, and Economies of Scale in Larger Multiplexers," *Queueing Systems*, vol. 20, pp. 293–320, 1995.

[14] G. L. Choudury, D.M. Lucantoni, and W. Whitt, "Squeezing the Most Out of ATM," *IEEE Transactions on Communications*, vol. 44, no. 2, pp. 203–217, 1996.

[15] N. G. Duffield, "Economies of Scale in Queues with Sources Having Power-Law Large Deviation Scaling," *Queueing Systems*, vol. 33, pp. 840–857, 1996.

[16] K.R. Krishnan, "A New Class of Performance Results for a Fractional Brownian Traffic Model," *Queueing Systems*, vol. 22, pp. 277–285, 1996.

[17] I. Saniee, A. Neidhardt, O. Narayan, and A. Erramilli, "Multiscaling models of sub-frame VBR video and TCP/IP traffic," *KICS/IEEE Journal of Communication Networks*, vol. 3, no. 4, pp. 383–395, 2001.

[18] M. Listani, V. Eramo, and R. Sabella, "Architectural and Technological Issues for Future Optical Internet Networks," *IEEE Communications Magazine*, vol. September, pp. 82–86, 2000.

[19] J. Cao, W. S. Cleveland, D. Lin, and D. X. Sun, "The Effect of Statistical Multiplexing on the Long-Range Dependence of Internet Packet Traffic," Tech. Rep., Bell Labs, Murray Hill, NJ, 2002.

[20] J. Cao, W. S. Cleveland, and D. X. Sun, "S-Net: A Software System for Analyzing Packet Header Databases," in *Proceedings Passive and Active Measurement*, 2002.

[21] Jin Cao and Kavita Ramanan, "A Poisson Limit for the Unfinished Work of Superposed Point Processes," in *Proceedings INFOCOMM*, 2002, to appear.

[22] W. S. Cleveland, D. Lin, and D. X. Sun, "IP Packet Generation: Statistical Models for TCP Start Times Based on Connection-Rate Superposition," *ACM SIGMETRICS*, pp. 166–177, 2000.

[23] D. R. Cox, *Renewal Theory*, Chapman and Hall, 1962.

[24] D. J. Daley and D. Vere-Jones, *An Introduction to the Theory of Point Processes*, Springer-Verlag, New York, 1988.

[25] J. Cao, W. S. Cleveland, and D. X. Sun, "Fractional Sum-Difference Models for Open-Loop Generation of Internet Packet Traffic," Tech. Rep., Bell Labs, Murray Hill, NJ, 2002.

[26] K. Claffy, H.-W. Braun, and G. Polyzos, "A Parameterizable Methodology for Internet Traffic Flow Profiling," *IEEE Journal on Selected Areas in Communications*, vol. 13, pp. 1481–1494, 1995.

[27] V. Paxson, "End-to-End Internet Packet Dynamics," in *Proceedings ACM SIGCOMM*, 1997, pp. 139–152.

[28] C. Fraleigh, C. Diot, B. Lyles, S. Moon, P. Owezarski, D. Papagiannaki, and F. Tobagi, "Design and Deployment of a Passive Monitoring Infrastructure," in *Proceedings Passive and Active Measurement*, Amsterdam, 2001, Ripe NCC.

[29] J. Micheel, S. Donnelly, and I. Graham, "Timestamping Network Packets," in *Proceedings ACM SIGCOMM Internet Measurement Workshop*, San Francisco, 2001.

[30] G.E.P. Box and G. M. Jenkins, *Time Series Analysis, Forecasting and Control*, Holden Day, San Francisco, 1970.

[31] J. R. M. Hosking, "Fractional Differencing," *Biometrika*, vol. 68, pp. 165–176, 1981.

[32] W. S. Cleveland and C. Liu, "Maximum Likelihood Estimation of Sum-Difference Time Series Models Using the EM Algorithm," Tech. Rep., Bell Labs, 2002.

[33] W. S. Cleveland, *Visualizing Data*, Hobart Press, Summit, New Jersey, U.S.A., 1993.