

Internet Traffic and Multiresolution Analysis

Ying Zhang¹, Zihui Ge², Suhas Diggavi³, Z. Morley Mao¹, Matthew Roughan⁴, Vinay Vaishampayan⁵, Walter Willinger⁵, and Yin Zhang⁶

Abstract: Traditional Internet traffic studies have primarily focused on the temporal characteristics of packet traces as observed on a single link within an ISP's network. They have contributed to advances in the areas of self-similar stochastic processes, long-range dependence, and heavy-tailed distributions and have demonstrated the benefits of applying a wavelet-based multiresolution analysis (MRA) approach when analyzing these traces. However, an ISP's physical infrastructure typically consists of 100's or 1000's of such links which are connected by routers or switches, and the Internet as a whole is made up of about 20,000 such ISPs. When viewed within this bigger context, the importance of the traffic's spatial characteristics becomes evident, and traffic matrices—compact and succinct descriptions of the traffic exchanges between nodes in a given network structure—are used in practice to capture and explore critical aspects of this spatial component of Internet traffic. In this paper, we first review some of the known results about the observed multifaceted scaling behavior of Internet traffic as seen on a single link. Next, we give a detailed account of how the architectural design of the Internet gives rise to natural representation of traffic matrices at different scales or levels of resolution. Moreover, we discuss the development of a MRA-like framework of traffic matrices that respects the different physically or logically meaningful Internet connectivity structures and provides new insights into Internet traffic as a spatio-temporal object.

Contents

1	Introduction	2
2	Single-link traffic: Self-similarity and Kurtz's construction	3
3	Network-wide traffic: Traffic matrices	7
	3.1 Traffic matrices at different levels of resolution	7
	3.2 Towards a MRA of traffic matrices	9
	3.3 A look at real traffic matrices	12
4	Summary and Outlook	15
	Acknowledgments	17
	References	17

¹Ying Zhang and Z. M. Mao are with the EECS Department, University of Michigan, Ann Arbor, MI 48109, USA, e-mail: {wingying,zmao}@eeecs.umich.edu

²Z. Ge is with Adverplex, Inc., Wakefield, MA 01880, USA, e-mail: {gezihui@adverplex.com}

³S. Diggavi is with the School of Computer and Communication Sciences, EPFL, Lausanne, SWITZERLAND, e-mail: suhas.diggavi@epfl.ch

⁴M. Roughan is with the School of Mathematical Sciences, University of Adelaide, Adelaide 5005, AUSTRALIA, e-mail: matthew.roughan@adelaide.edu.au

⁵V. Vaishampayan and W. Willinger are with AT&T Labs-Research, Florham Park, NJ 07932, USA, e-mail: {vinay,walter}@research.att.com

⁶Yin Zhang is with the Department of Computer Sciences, University of Texas at Austin, Austin, TX 78712, USA, e-mail: yzhang@cs.utexas.edu

AMS 2000 subject classifications: Primary 60K30, 90B18; secondary 60G18, 60G57

Keywords and phrases: Self-similar processes, heavy-tailed distributions, traffic matrices, wavelet-based multi-resolution analysis

1. Introduction

Internet traffic is a multi-faceted object, and depending on one's vantage point, can either be viewed as a purely temporal, a purely spatial, or a combined temporal-spatial process. Traditional Internet traffic studies have focused mainly on the traffic that traverses a given link between two routers in the network. Their primary objective has been to describe the pertinent statistical characteristics of the temporal behavior of the measured traffic rate process (i.e., the number of bytes or packets seen on the given link in successive time intervals). However, Internet traffic arises in a very structured manner that reflects the architectural design of the Internet, with the vertical separation into layers and the horizontal decentralization across network components representing two of its most prominent and influential features. As a result, the focus of subsequent Internet traffic studies has turned from being largely descriptive to being fully explanatory in the sense that observed characteristics of traffic rate processes are traced to particular aspects or mechanisms that determine how traffic is generated and handled within the confines of the architectural framework provided by today's Internet.

As an example, one of the most visible manifestations of the Internet's vertical decomposition is the 5-layer TCP/IP protocol stack, consisting of (from the bottom up) the physical layer (e.g., optical fiber, copper), the data link or network access layer (e.g., Ethernet, frame relay), network or internet layer (e.g., *Internet Protocol*, or IP), the host-to-host or transport layer (e.g., *Transmission Control Protocol*, or TCP), and the application layer (e.g., *HyperText Transfer Protocol*, or HTTP, for the World Wide Web, or WWW). In turn, on a given link, every bit recorded at the physical layer can in general be uniquely associated with higher-layer entities such as IP packets, IP flows, TCP connections, or application-layer sessions. As discussed in Section 2, the challenge of explanatory Internet traffic modeling has been to relate observed features of measured traffic rate processes to and explain them in terms of properties of these higher-layer traffic entities. Key to these efforts has been *Kurtz's construction* [31], a flexible framework for exploring Internet traffic that (i) is mathematically rigorous, (ii) accounts for the different layers in the TCP/IP protocol stack, (iii) is consistent with measured Internet traffic across the different layers, and (iv) highlights the intimate connection between the observed temporal scaling properties of various traffic rate processes and the ubiquitous high-variability or heavy-tailed properties of the different higher-layer entities (e.g., IP flows, TCP connections, sessions). These efforts have been aided by the development and application of a *1D wavelet-based multi-resolution analysis (MRA)* that has enabled an in-depth examination of the temporal dynamics of measured Internet traffic across a wide range of time scales of interest [2]. We summarize the main findings and implications from these efforts in Section 2.

The main objective of this paper is to outline a similar wavelet-based MRA (in 2D instead of 1D) for studying the spatial features of a network's measured traffic matrices rather than the temporal aspects of the measured traffic rate processes on a link. Here, a traffic matrix describes the amount of traffic (in bytes or packets) that is sent from one point or node in a network to another during some time interval (e.g., 5-30 minutes). The networks of interest are manifestations of the Internet's horizontal decentralization and reflect physically, logically, or managerially meaningful ways of organizing the Internet-wide physical infrastructure into smaller entities or subnets. In turn, nodes can represent physical links, routers, Points-of-Presence (PoPs), autonomous systems (ASs) or domains, or entire ISPs, and the resulting traffic matrices provide compact descriptions of network-wide In-

Internet traffic across a wide range of spatial scales of interest. In Section 3, we discuss traffic matrices at different levels of scale or resolution, outline a rudimentary *2D wavelet-based MRA* of traffic matrices, and illustrate it with some examples of actual traffic matrices. In particular, we use the Abilene network [1] since Abilene makes all of the information about its network publicly available. However, we hasten to point out that the development of a MRA of traffic matrices will benefit from future studies that examine a wider range of networks, especially from commercial ISPs. In addition, many technical and practical problems remain, and we conclude in Section 4 with a list of the most pressing open problems and a discussion of promising applications of the envisioned MRA of traffic matrices.

2. Single-link traffic: Self-similarity and Kurtz's construction

When focusing on a single link within an ISP's network, Internet traffic data typically consists of high time-resolution measurements recorded on that physical link over which the bits are sent. Each transmitted bit seen on this link can in general be associated with higher-layer entities such as an IP packet. Along with every packet header that is captured and stored, additional information is usually saved, notably an accurate time stamp (packet arrival time), packet size (number of bits or bytes), and other status and possibly even some payload information. Empirical studies of these high-quality and high-volume data sets have generally focused on identifying and describing pertinent statistical characteristics of the temporal dynamics of the measured packet or bit rate processes (i.e., the time series representing the number of packets or bits per time unit, over a certain time interval). They have provided ample evidence that measured Internet traffic exhibits extended temporal correlations (i.e., *long-range dependence*), and that when viewed within some range of moderately small to moderately large time scales, the traffic appears to be *fractal-like* or *self-similar*, in the sense that a segment of the traffic measured at some time scale looks or behaves just like an appropriately scaled version of the traffic measured over a different time scale.

The original finding of self-similar scaling behavior in measured network traffic was reported in [24, 32] and was based on an extensive statistical analysis of traffic measurements from Ethernet local-area networks over a four-year period from 1989-1993 [32, 33]. A number of key follow-up studies have provided further evidence of the prevalence of self-similar traffic patterns in measured pre-Web Internet traffic [42, 43] and post-1995, Web-dominated Internet traffic [11, 12] (see also [54, 41, 15] and references therein) and have contributed to a general acceptance of self-similarity as an *invariant* [23] of Internet traffic—a traffic characteristic that has been largely insensitive to the sometimes drastic changes the network and its traffic have undergone during the past 10 or so years. Subsequent empirical studies have refined this picture by focusing on measured network traffic over very large as well as over very small time scales. Over large time scales (e.g., hours or days), traffic has been found to be largely dominated by pronounced time-of-day and day-of-week effects (e.g., see [46, 48]), and traffic models used for network engineering purposes such as link dimensioning and capacity planning need to account for this property of network traffic. With respect to the dynamic nature of network traffic over small time scales (i.e., below the typical round-trip time (RTT) of a packet), recent work has demonstrated that it also deviates from the self-similar scaling behavior that has been observed over larger than RTT time scales. However, in contrast to incorrect claims about an apparent Poisson-like dynamics of Internet

traffic on sub-RTT time scales (e.g., see [9, 30]), there generally exists significant burstiness even on very small time scales, mainly due to the closed-loop feedback dynamics inherent in TCP-dominated Internet traffic (e.g., [19, 21, 22, 17]) and the intricate traffic interactions that can occur within a router and across the network. On the one hand, this understanding has provided new insight into when the use of self-similar traffic models such as fractional Brownian motion (or equivalently, its increment process, fractional Gaussian noise) is justified and can be exploited for network engineering purposes [18]. On the other hand, it has also highlighted the need for a paradigm shift in network traffic modeling, whereby the currently employed strictly open-loop traffic models need to be replaced by constructs that can account for the critical features of TCP-type feedback regulation. The importance of considering relevant closed-loop traffic models becomes apparent when studying networking problems such as adequate sizing of router buffers in today’s networks (e.g., see [27, 7, 17] or helping a service provider to offer and guarantee competitive service-level agreements (SLAs) to its customers [50].

In effect, these empirically-based efforts toward describing actual Internet traffic have demonstrated that self-similar processes, despite their limitations on both sides of the spectrum of relevant time scales, define an elegant family of compact mathematical models for capturing the essence behind the wide range of “burstiness” or scale-invariance encountered in measured traffic traces. In turn, the self-similarity discovery has invigorated research in the area of statistics for long-range dependent and self-similar stochastic processes (e.g., see [8, 49]). More importantly, however, it has motivated the construction of new mathematical models that provide a physical (i.e., networking-based) explanation of the observed self-similar scaling behavior of Internet traffic that is intuitively appealing, conceptually simple, mathematically rigorous, and verifiable. Recognizing that it is difficult to think of many other areas in the sciences where the available data provides such detailed information about so many different facets of the system under study, these models have by and large succeeded in demystifying self-similarity as an Internet traffic characteristic by explicitly accounting for key aspects of the design and architectural principles of today’s Internet and enabling direct model validation that relies on and exploits the high semantic context contained in the measured data. They are in stark contrast to the traditional traffic models that are “black boxes” in the sense that they ignore nearly all of this rich semantic context (they tend to use only packet arrival time and packet size information), describe the traffic traces at hand well in a statistical sense, but typically contribute little or nothing to our understanding of data networks and the traffic they carry.

To illustrate, given accurate packet header information, measured Internet traffic can be sliced and diced in many different ways, resulting in a number of different representations of network traffic as seen on a single link. For example, by extracting packet header-specific information such as source- and/or destination IP address or prefix, source- and/or destination port numbers, protocol-, or application-specific attributes, it is possible to uniquely associate each packet with the IP flow, TCP connection, and sometimes even application-layer entity that it belongs to. Such mappings decompose traffic naturally into individual constituents that have network-specific meaning and explicitly reflect the layered design of the Internet architecture—IP flows allow for an IP layer view of traffic, TCP connections for a TCP layer perspective, and higher-layer entities provide a glimpse into network traffic at the application layer. In turn, such decompositions support model formulations that treat traffic on a link as an aggregate of many such constituents and invite model constructions that don’t view packets as black boxes but are given

in terms of these constituents. An elegant mathematical framework for generating and analyzing such physical or “structural” models is due to Kurtz [31]. *Kurtz’s construction* considers traffic models that are integral representations with respect to certain Poisson random measures, and in its general form, includes well-known earlier approaches such as *Cox’s construction* [10] (also known as immigration death process or $M/G/\infty$ queuing model) and *Mandelbrot’s construction* [37] (also known as renewal-reward process).

In its basic form, Kurtz’s construction accounts for the layering architecture of the Internet by assuming for example that at the application layer, sources or *sessions* (e.g., FTP, HTTP, TELNET) arrive at random (i.e., according to some stochastic process) on the link and have a “lifetime” or session length during which they exchange information. At the IP layer, this information exchange manifests itself as a flow of IP packets that are transmitted at, say, some constant rate from the start until the end of a session. Thus, at the IP layer, the aggregate link traffic measured over some time period is made up of the contributions of all the sources that—during the period of interest—actively transmitted packets. More formally, one representation of this aggregate link traffic or workload process is as follows. Let the source activation process be $N = (N(t) : t \geq 0)$, denoting the number of source activations up to time t ; for the i -th activation, let $X_i(s)$ denote the total traffic generated by source i during the first s units of time. We model the length of time τ_i that source i remains active separately from X_i and assume that the pairs (X_i, τ_i) are *i.i.d.* The total link traffic or workload generated up to time t can then be written as

$$(2.1) \quad U(t) = \int_0^t X_{N(s)}(\tau_{N(s)} \wedge (t-s)) dN(s),$$

and if $L = L(t) : t \geq 0$ denotes the number of sources that are active at time t , we have

$$(2.2) \quad L(t) = \int_0^t I_{[0, \tau_{N(s)}]}(t-s) dN(s).$$

Assuming that N is a counting process with intensity $\lambda(N, U, L, \cdot)$, that is,

$$(2.3) \quad N(t) - \int_0^t \lambda(N, U, L, s) ds$$

is a martingale with respect to the filtration generated by the random variable $\{N(s), U(s), L(s); s \leq t\}$, the process (N, U, L) can be represented as the solution of a system of stochastic equations involving a Poisson random measure (see [31] for details).

This representation provides a convenient mathematical framework for studying scaling limits of the process (N, U, L) that yield deterministic “fluid approximations” or corresponding central limit theorems. For example, by appropriately scaling session intensity and time, the workload process U can be shown to converge to a self-similar limiting process, namely *fractional Brownian motion* (or its increment process, *fractional Gaussian noise*), provided the session arrivals follow a Poisson process, the sessions share the bandwidth in a “fair” (i.e., TCP-like) manner, and, more importantly, the session durations or lifetimes are *i.i.d.* and have a distribution that is *heavy-tailed with infinite variance* [5, 45]. Intuitively, the latter condition implies that there is no “typical” session size but instead, the session sizes are “highly variable” (i.e., exhibit infinite variance) and fluctuate over a wide

range of scales, from Kilobytes to Megabytes and Gigabytes and beyond. It is this basic characteristic at the higher layers in the TCP/IP protocol stack that causes the aggregate traffic at the IP layer to exhibit self-similar scaling. By relaxing the fair bandwidth sharing assumption, allowing for more realistic within-session traffic rates, or manipulating the relative speed with which the number of sessions and the time scale increase, other types of (self-similar and not self-similar) limiting workload processes are possible, including *Levy-stable motion* or its increment processes, *fractional Levy-stable noise* (for more details, see [31, 34, 52, 39, 44, 28]).

The beauty of structural models such as Kurtz's construction is that in stark contrast to the conventional black box models, they not only explain the self-similarity phenomenon in simple terms (i.e., heavy-tailed connections), but they also clearly identify the data sets that need to be either obtained from new measurements or extracted from the available IP packet-header traces to validate the proposed explanation. This "closes the loop" between the empirical discovery of the self-similar scaling behavior of aggregate Internet traffic on the one hand, and its mathematical explanation in terms of infinite variance phenomena associated with meaningful quantities at the higher layers in the TCP/IP protocol stack on the other. For example, because of the way many applications are structured, determining session-related entities such as arrival times and sizes or durations from packet-level measurements is straightforward. For FTP and TELNET, these entities have been shown to be consistent with Kurtz's construction in [43]. For HTTP (i.e., Web sessions), obtaining session information is generally more involved [29], but the empirical evidence for the heavy-tailed characteristic of Web-related entities (e.g., HTTP request sizes and durations) has been well-established to date (see for example [56, 12, 14, 57]). In fact, heavy-tailed characteristics of higher-layer entities such as IP flows, TCP connections, or sessions constitute yet another set of Internet traffic invariants.

While the self-similar scaling behavior across a range of intermediate time scales of Internet traffic at the IP layer (i.e., the time series representing the number of packets or bytes per time unit) is well documented, an equally intriguing scaling property of Internet traffic across the higher layers in the TCP/IP protocol stack has received comparatively little attention. For example, instead of viewing Internet traffic at the IP layer in terms of a time series representing the number of IP packets per time unit, we can consider physically meaningful "coarsened" versions by, for example, defining Internet traffic at the IP layer as given by the time series representing the number of IP flow arrivals per time unit or, for that case, Internet traffic at the TCP layer as given by the time series representing the number of TCP connection arrivals per time unit. As originally pointed out in [20], the latter also exhibit self-similar scaling characteristics which have in fact become more pronounced as the traffic mix at the application layer has changed from mostly TELNET and simple use of email and FTP during the pre-Web period to predominantly Web-based after about 1995 [20, 21]. Note that Kurtz's construction applies equally well for explaining the self-similar scaling behavior of these coarsened versions of Internet traffic as observed at the IP and TCP layers, respectively, and simply requires the distribution of the number of IP flows (or TCP connections) per session to be heavy-tailed with infinite variance [21].

Together, these observations suggest that the different self-similar scaling phenomena observed in measured Internet traffic are mainly caused by user/application characteristics, have little to do with the network (except that it imposes some fair sharing of bandwidth), and are likely to remain with us in the foreseeable future. Note that to arrive at this basic understanding of the temporal dynamics of Internet

traffic as seen on a single link within the network, the development and application of a *1D wavelet-based MRA* in support of a detailed examination of measured Internet traffic over a wide range of time scales of interest has been of critical importance [2, 3, 4]. Acting as an *analytic telescope*, this wavelet-based MRA has been ideal for the study of scaling properties and as such, has enabled a data analysis that matches well with the properties encountered in measured Internet traffic. Moreover, this technique has been instrumental in demonstrating that alternative models that are capable of reproducing long-range dependencies or self-similar scaling behavior (e.g., conventional stochastic processes with built-in non-stationarities such as deterministic monotonic trends or level shifts of the mean) are by and large inconsistent with measured network traffic (e.g., see for example [55, 2]). This development has been accompanied by equally important advances in the area of inference for heavy-tailed phenomena (e.g., [45, 5]). In particular, the high-volume of the available data sets has motivated a pragmatic approach to dealing with high-variability in network measurements that has its roots in Mandelbrot's early work [36]. This approach is described in [58] and clarifies in which sense higher layer traffic quantities such as sizes or durations of IP flows, TCP connections, or sessions are fully consistent with proper infinite variance distributions, but are by and large inconsistent with conventional, finite variance distributions such as Lognormal or Weibull distributions. Thus, as far as the self-similar scaling behavior of Internet traffic is concerned, the explanation in terms of high-variability phenomena (i.e., infinite variance distributions) at the higher layers in the protocol stack remains to date the only model that is mathematically rigorous and consistent with measured network traffic as it manifests itself (in different incarnations) across the different layers of the protocol stack.

3. Network-wide traffic: Traffic matrices

The traffic observed on a single link within an ISP's network arrives at that link coming from possibly many different sources and leaves the link destined to possibly many different destinations. While our main focus in Section 2 was on the temporal dynamics of Internet traffic as observed on a single link, here we are less concerned with its temporal aspects, but are mainly interested in its spatial properties. To this end, the objects of interest are traffic matrices [38, 6], and to simplify the presentation, we discuss traffic matrices in the context of a single network or domain (i.e., a single network that connects end-systems, but does not connect to other networks).

3.1. Traffic matrices at different levels of resolution

Traffic matrices describe the amount of traffic from one point in a network to another during some time interval, and are thus naturally represented by a three-dimensional data structure $T_t(i, j)$ which represents the traffic volume (in bytes or packets) from i to j during a time interval $[t, t + \Delta t)$. The locations i and j are generally considered to be discrete in nature, i.e. they are drawn from some set of possible locations. We may consider these locations to be physical geographic locations, thereby making i and j spatial variables. However, in the Internet, it is common to associate i and j with logical structures related to the address structure of the Internet, i.e. IP addresses, or natural groupings of such by a common prefix corresponding to a subnet.

In the case that locations are modeled spatially, we call the set of possible locations L , while the set of address-based locations will be denoted I . In general, there is some correspondence between the two, but the mapping is not one-to-one, and so we cannot in general map a location to an IP address, or vice versa.

Given some set of locations I , we can easily aggregate a traffic matrix across sets $S, D \subset I$ to obtain

$$T_t(S, D) = \sum_{i \in S} \sum_{j \in D} T_t(i, j).$$

As with any time series, we can also perform standard aggregation operations in time, making it relatively straightforward to create multiple approximate representations of the original traffic matrix at different levels of resolution.

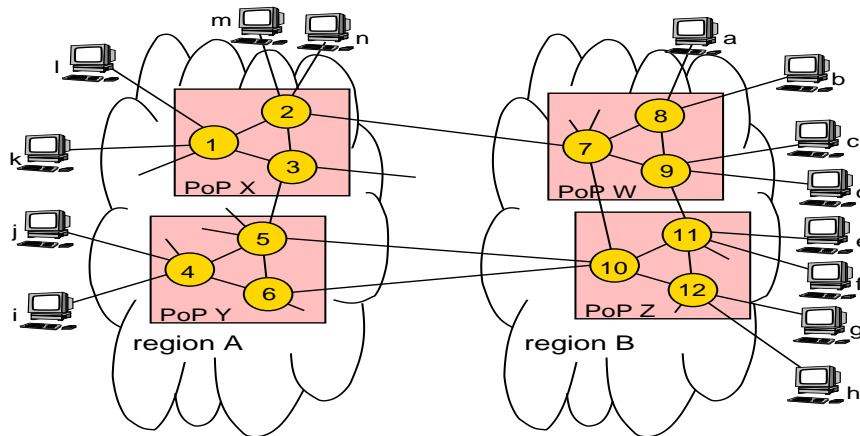


FIG 1. Example network.

However, such an approach, while potentially useful, might provide only limited additional insight into the nature of traffic matrices. The key to interesting approximation lies in the choice of sets S, D used at each step in the aggregation of the matrices. The reason for this lies in the designed structure of a network. For instance, consider the network in Figure 1. The figure shows a toy network comprising two regional networks, where each subnet contains several Points-of-Presences (PoPs), each of which in turn contains a number of routers, which connect to multiple end-systems. It seems obvious that this purposefully engineered hierarchy should be related to the manner in which we perform the “coarsening” of traffic matrices.

For instance, in Figure 1, we might naturally consider the end systems to be the locations of interest, i.e. $A = \{a, b, c, d, e, f, g, h, i, j, k, l, m\}$, and then aggregate first by router, so that we take sets

$$\begin{aligned} S_1^{(1)} &= \{l, k\}, S_2^{(1)} = \{m, n\}, S_3^{(1)} = \{\}, \\ S_4^{(1)} &= \{i, j\}, S_5^{(1)} = \{\}, S_6^{(1)} = \{\}, \\ S_7^{(1)} &= \{\}, S_8^{(1)} = \{a, b\}, S_9^{(1)} = \{c, d\}, \\ S_{10}^{(1)} &= \{\}, S_{11}^{(1)} = \{e, f\}, S_{12}^{(1)} = \{g, h\}. \end{aligned}$$

Note that, in reality there would likely be many more end-systems, and hence the sets would be rather larger. We could then aggregate into PoPs, such that

$$S_X^{(2)} = S_1^{(1)} \cup S_2^{(1)} \cup S_3^{(1)},$$

and regions such that

$$S_A^{(3)} = S_X^{(2)} \cup S_Y^{(2)},$$

and so on for other regions and PoPs.

Note that the superscript in the sets above is used to denote the level (scale) of resolution (approximation) that would be involved in calculating $T_t(S_a^{(i)}, S_b^{(i)})$. We will typically denote a traffic matrix aggregated across sets $S_a^{(i)}$ by $T^{(i)}$, but we will also retain this notation for any approximation at level i .

In the network above, the topological hierarchy is defined by the administrator and has some meaning, either geographically, or managerially. However, it may not be obvious in some networks what the natural groupings are. For instance, regions may not be well defined in many networks. In this case, it may make sense to group the end hosts using a clustering algorithm based on network distances: many networks use shortest-path routing where the link weights (distances) are administratively defined, and these distances between end-points define a natural clustering, or hierarchy on the network. Similarly, there may be circumstances where the logical hierarchy is more important than the physical topology when creating approximate representations of the network. For example, IP addresses have a natural hierarchy, which does not necessarily mesh with geography. Another example occurs when end-points connect to multiple points in the network (for redundancy), and it might make sense to aggregate over these logically related end-points. In particular cases, we may be able to define a natural logical hierarchy suitable for generating appropriate approximations. However, in many cases, there may be no obvious logical hierarchy, in which case, part of our goal may be to search for the “best” hierarchical decomposition.

3.2. Towards a MRA of traffic matrices

Multi-Resolution Analysis (MRA) (or alternatively Multi-Resolution Approximation) refers to the process of creating multiple approximate representations of an object (e.g., traffic matrix), such that these have different resolution. In the well-known context of wavelets, fast algorithms exist to calculate these approximations at a countable number of resolutions (for a different approach using wavelets for spatial traffic analysis, see [13]). MRA can be useful for a number of problems, including denoising, compression, and anomaly detection. Here we wish to extend these ideas to Internet traffic matrices so that we may be able to use these types of applications in practice, but more importantly, to gain fundamental insight into the nature of actual Internet traffic matrices.

However, there is more to MRA than simple aggregation/approximation, and one of the main objectives is to be able to decompose a given traffic matrix in such a way that when we form our successive approximate representations (using a “decomposition algorithm”), we can also retain enough information to reverse the approximation process—in wavelet parlance, we wish to retain the *details* necessary for obtaining high-fidelity reconstructions (using a “reconstruction algorithm”) [16, 53]. In essence, the approach is intended to find sparse representations of traffic matrices such that one can represent their important features with a small set of numbers. Our objective is to understand traffic matrices at a level which will aid in synthesis (artificial generation of traffic matrices for the purpose of simulations [47]) or inference (statistical estimation of a traffic matrix from link load data [59]). In this context, sparse representations of traffic matrices are of special interest. For

example, for synthesis, they reduce the number of parameters that must be tuned or estimated. For inference, the key problem is the massively under-constrained nature of the linear-inverse problem that must be solved—if we can reduce the problem to inference of a smaller number of parameters, then it will no longer be under-constrained.

To illustrate some of the features of the envisioned MRA of traffic matrices, we are motivated by existing work in image compression. To this end, consider a $n \times n$ traffic matrix T and let G represent the *analysis matrix* of a wavelet transform. For example, when using the Haar transform, the analysis matrix is given by

$$(3.1) \quad G = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 & 0 & \dots & 0 & & & & & & \\ 0 & 0 & 1 & 1 & 0 & \dots & \dots & & & & 0 \\ \vdots & & & & & & & & & \vdots & \vdots \\ 0 & 0 & & & 0 & & & & & 1 & 1 \\ 1 & -1 & 0 & \dots & 0 & & & & & & \\ 0 & 0 & 1 & -1 & 0 & \dots & \dots & & & & 0 \\ \vdots & & & & & & & & & & 1 & -1 \end{pmatrix}.$$

While generalizations of the wavelet transform to multiple dimensions are known, one of the simplest methods for applying such a transform in higher dimensions is in a separable fashion. That is, given a traffic matrix T , a separable wavelet transform matrix is computed by applying the analysis matrix G to the rows and columns of T separately. This results in the matrix A given by

$$(3.2) \quad A = GTG^t,$$

where G^t is the transpose of G . One of the reasons why wavelet transforms are widely used is that the resulting wavelet representations tend to concentrate energy in a few of the coefficients, thus resulting in a sparse representation. To explain this in more detail, consider a partition of the matrix G into two block submatrices of size $n/2 \times n$; that is,

$$(3.3) \quad G = \begin{pmatrix} G_1 \\ G_2 \end{pmatrix}.$$

This in turn allows us to partition the wavelet transformed matrix A as follows

$$(3.4) \quad A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$$

where $A_{ij} = G_i T G_j^t$. The wavelet coefficients in each of the four submatrices allow for different interpretations, exhibit different behavior, and can be quantized differently in order to gain a compression advantage. For example, as in image compression, where using A_{11} alone to reconstruct the image will produce a smoothed version of the original image, in the case of traffic matrices, A_{11} defines an approximate ‘‘coarsened’’ version of the original traffic matrix that is obtained by aggregating over appropriate rows and columns of T . On the other hand, the details are contained in the other submatrices, and their contributions to the reconstructed traffic matrix can be controlled by thresholding them and by tuning the value of the threshold. In the area of image compression, much is known about the subtle correlations between the wavelet coefficients, especially when the wavelet transform

is applied repeatedly. However, in the context of traffic matrices, such an understanding is still missing.

Note that separability in the context of traffic matrices has a clear interpretation. Aggregating rows corresponds to aggregating source nodes and aggregating columns corresponds to aggregating destination nodes. In effect, separability allows for a 2D transform that reduces to two 1D transforms, one for source nodes, the other for destination nodes. Since the rows and columns selected for aggregation by the analysis matrix G are fixed and may not correspond to any physically, managerially, or logically defined hierarchy, we can provide some extra flexibility by permuting the rows and columns using permutation matrices Π_r and Π_c^t , respectively. This results in more flexible wavelet transformed matrices of the form

$$(3.5) \quad A = G\Pi_r T \Pi_c^t G^t.$$

This of course leads to the problem of permutation matrix selection, a hard combinatorial problem by itself, unless the permutations to use are obvious and arise naturally within the hierarchical network structure of interest.

As far as reconstruction is concerned, the wavelet transformed matrix A may be inverted through the use of the synthesis matrix $H = G^{-1}$, via the equation $T = HAH^t$. The wavelet transform is information preserving and thus T can be recovered exactly¹. However, from a modeling perspective, we would like to study the quality of the resulting traffic matrix estimate when the wavelet coefficients are modified in some way, e.g., by setting some to zero through an appropriately designed thresholding operation. This requires that adequate distance measures be used (e.g., Kullback-Leibler divergence, norm-based metrics such as l_2 -norm or Frobenius norm), but a detailed study of appropriate evaluation metrics is beyond the scope of this paper and will appear elsewhere. As an example, in Section 3.3 below, we will consider a reconstruction of the form

$$(3.6) \quad \hat{T} = H \begin{pmatrix} A_{11} & 0 \\ 0 & 0 \end{pmatrix} H^t.$$

Concerning the structure of the wavelet transformed traffic matrix itself, note that in the case of the Haar analysis matrix, A_{11} is guaranteed to have non-negative entries (since the entries in G_1 are non-negative) and can therefore be thought of as a genuine traffic matrix. The other submatrices carry detail information that is lost in the aggregation and will in general have non-negative as well as negative entries.

A traffic matrix T is called a *gravity matrix* or *gravity model* if $T = uv^t$ for some vectors u and v . Gravity models have been used successfully as models for traffic matrices in inference and synthesis [47], though they have limitations [6]. Note that the defining property of a gravity matrix corresponds to separability of the traffic matrix. Moreover, if T is a gravity matrix, we have

$$(3.7) \quad A = (Gu)(Gv)^t,$$

that is, A is also a gravity matrix. Since this is independent of the precise type of analysis matrix G , the property of being a gravity matrix is preserved not only under aggregation, but under other forms of filtering as well. Similarly, it also follows

¹This is assuming infinite precision arithmetic if the wavelet analysis matrix contains irrational numbers. The desire for exact reconstruction is one of the motivations for considering *lifting* schemes.

that each of the submatrices A_{ij} is also a gravity matrix. In terms of reconstructed traffic matrices, it is not hard to see that the reconstruction (3.6) will result in a gravity matrix provided the submatrix matrix A_{11} is a gravity matrix. On the other hand, the reconstruction

$$(3.8) \quad \hat{T} = H \begin{pmatrix} A_{11} & 0 \\ 0 & A_{22} \end{pmatrix} H^t.$$

will not be a gravity matrix if A_{22} is nonzero. Thus if we wish the approximate traffic matrix

$$(3.9) \quad \hat{T} = H \begin{pmatrix} \tilde{A}_{11} & \tilde{A}_{12} \\ \tilde{A}_{21} & \tilde{A}_{22} \end{pmatrix} H^t$$

to be a gravity matrix, we must be careful to ensure that \tilde{A} (with submatrices \tilde{A}_{ij}) is also a gravity matrix.

Gravity models are of particular interest in the context of MRA of traffic matrices because the underlying assumption of the gravity model (i.e., traffic homogeneity) is expected to improve with aggregation. Larger aggregates of traffic should behave more and more like a gravity model, until the top level approximation (just the total traffic in the network) is exactly represented by such a model. Note however that systematic biases away from a gravity model may be regional, so aggregating topologically may actually result in delayed convergence to the gravity model, whereas, randomized aggregation may actually converge quite quickly to fit a gravity model. Other approaches to aggregation that are less oblivious to actual routing of the traffic through the network may have the benefit of quick convergence **and** lack of systematic bias.

3.3. A look at real traffic matrices

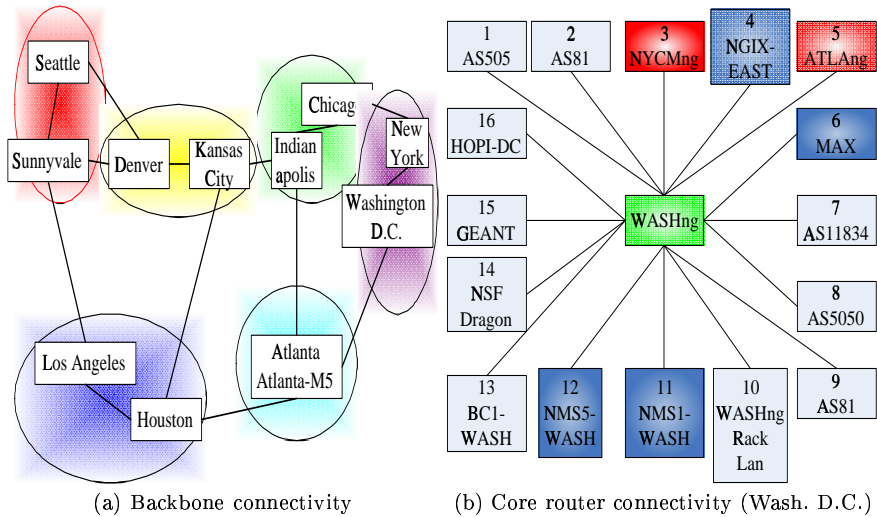


FIG 2. The Abilene network.

To illustrate various features of actual traffic matrices, we first consider the Abilene network shown in Figure 2. Abilene [1] is the U.S. Internet backbone for higher

education. It is comprised of high-speed connections between core routers (Juniper T640) which are located in 11 U.S. cities, with the Atlanta node consisting of two core routers (i.e., "Atlanta" and "Atlanta-M5"). The Abilene backbone is shown in Figure 2(a) and is a sparsely connected mesh; connectivity to regional and local customers is not shown but is provided with some minimal amount of redundancy. Abilene maintains peering connections with other higher educational networks (domestic and international) but does not connect directly to the commercial Internet. This feature is shown in Figure 2(b) which depicts the connectivity of Abilene's core router in Washington, D.C. at the level of populated router interfaces (numbered 1–16). For example, interfaces 3 and 5 connect to the core router in New York and one of the core routers in Atlanta, respectively; interfaces 4, 6, 11, and 12 connect to Internet exchange points; and the other interfaces shown reflect peering connections to customers such as AS81 which belongs to the North Carolina Research and Education Network NCREN-MCNC.

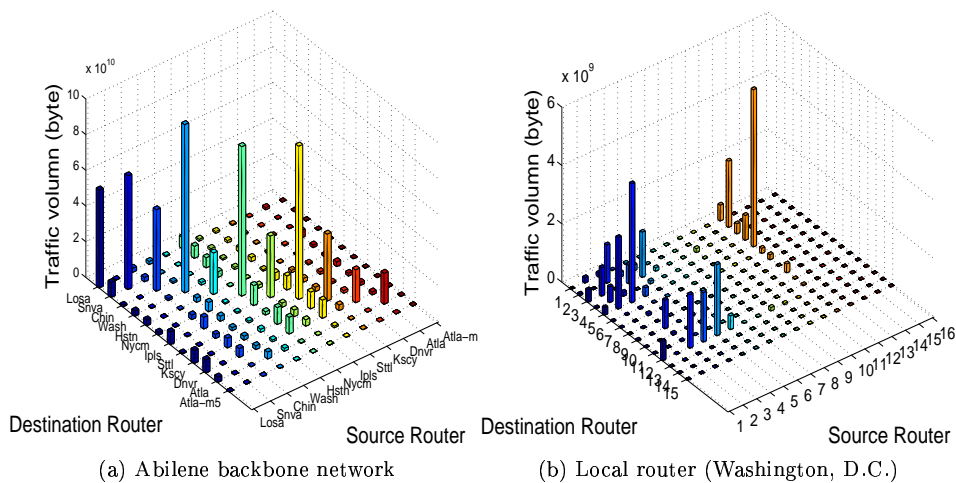


FIG 3. Measured traffic matrices.

A snapshot of Abilene's traffic matrix is shown in Figure 3(a) and depicts the amount of traffic carried between each Abilene node on 09/01/2006. For that same day, the local router traffic matrix for the Washington, D.C. node is shown in Figure 3(b). Note that the large diagonal elements in Figure 3(a) reflect a pronounced locality property of Abilene traffic, while the local router traffic matrix in (b) is largely determined by the configuration of this core router (i.e., which interface carries which in- and out-going traffic). Plotting in Figure 4(a) the values of the 12 largest elements of the traffic matrix in Figure 3(a) for successive 1-hour intervals for the 6-day period from 09/01/2006 to 09/06/2006 shows the presence of a dominant diurnal cycle that has been well-documented in past studies of single-link traffic dynamics over large time scales [48]. A very similar behavior can be observed for the entries on the local router traffic matrix in Figure 4(b).

Considering the static Abilene traffic matrix T in Figure 3(a), two simple approximations are obtained by computing the corresponding gravity model T_G and deriving the wavelet transformed model T_W of the form given by (3.6). Note that for the gravity model, we have $T_G = uv^t$ with $u_i = (1/\sqrt{(S)}) \sum_j T_{i,j}$ and $v_j = (1/\sqrt{(S)}) \sum_i T_{i,j}$, where $S = \sum_{i,j} T_{i,j}$. To derive a simple, yet meaningful

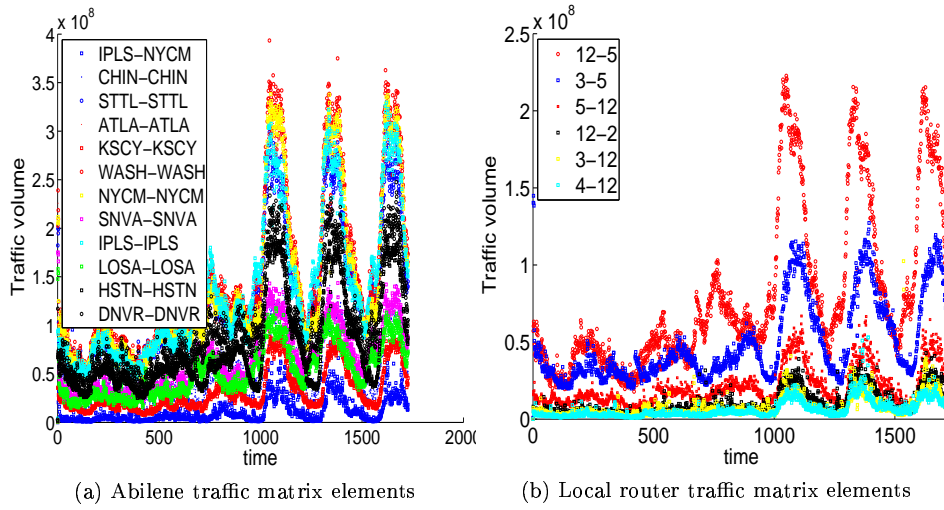


FIG 4. Measured traffic matrices over time.

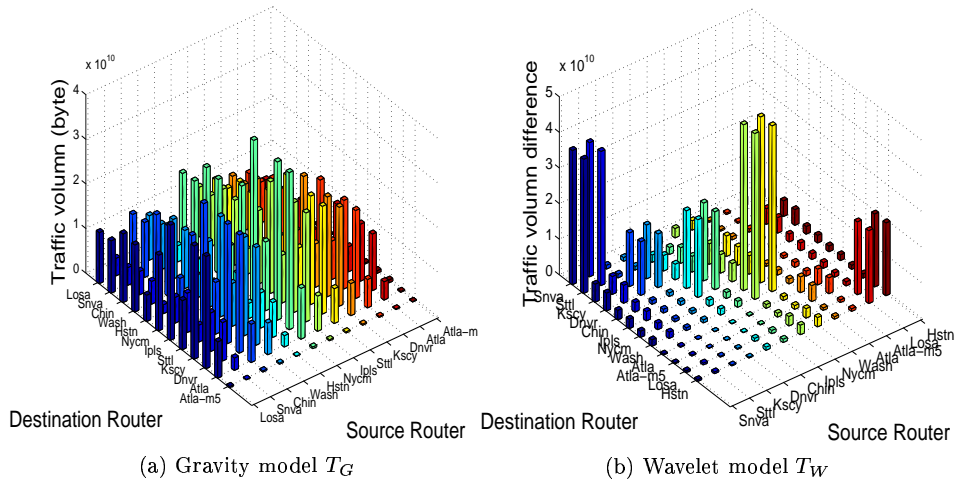


FIG 5. Approximate traffic matrices.

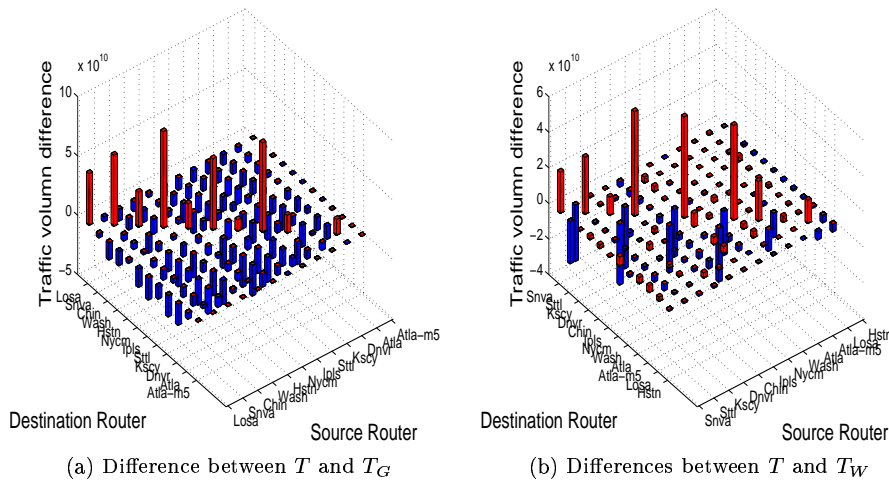


FIG 6. Quality of traffic matrix approximations; red bars (pointing up) are positive, blue bars (pointing down) are negative differences.

wavelet transformed traffic matrix, we aggregate the Abilene nodes geographically in pairs of two as shown in Figure 2(a) by using appropriate permutation matrices $\Pi_r = \Pi_c$, compute the wavelet transformed matrix A via (3.5), and set $T_W = \hat{T}$ where \hat{T} is given by equation (3.6). While Figure 5 shows the two approximate traffic matrices, Figure 6 depicts the differences between T and T_G , and between T and T_W , respectively. Note that while neither approximations can account for the large diagonal elements of T , the wavelet transformed traffic matrix T_W results in a qualitatively better approximation of T than the gravity model T_G . At the same time, Figure 6(b) also shows the effects of relying on the simple dyadic structure associated with the Haar transform when choosing the matrix A given by (3.1) as our analysis matrix. When comparing Figures 3(a) and 5(b), this aggregation into groups of two appears as the most significant difference between the original and the wavelet transformed traffic matrices and suggests alternate and more flexible choices of wavelet transforms and corresponding analysis matrices. However, not every choice that is meaningful from a networking perspective is feasible from an MRA perspective (i.e., A may not be invertible, causing problems for the reconstruction), and herein lies much of the tension that exists between developing a MRA that is, on the one hand, suitable for the Internet context and, on the other hand, amenable to a rigorous mathematical treatment.

4. Summary and Outlook

By combining the analysis of single-link traffic rate processes with the more recent studies of network-wide traffic matrices, a detailed exploration of Internet traffic as a spatial-temporal object across the different layers of the TCP/IP protocol stack looms as a real possibility. However, to study Internet traffic over a wide range of scales in space and time and across different layers will require a dramatic widening of MRA technology as it is known and used today. In Section 3, we discussed some basic features of such an MRA for the case of static traffic matrices, but much work remains even in this case where temporal and layer-specific aspects are largely suppressed and the focus is on the spatial characteristics of the total traffic volumes

exchanged between pairs of nodes in the network. In particular, we would like to know how to coarsify traffic matrices in such a way that the reconstructed approximations automatically satisfy the non-negativity constraints and can therefore be interpreted as genuine traffic matrices. In the case of wavelet transformed matrices, we are especially interested in thresholding techniques that ensure non-negativity of the reconstruction, especially when the transform is applied iteratively. Other open issues concern the choice of appropriate metrics for comparing different traffic matrices across scales and within a given scale; the development of flexible “zoom-in” capabilities for exploring Internet traffic localized in time, space, or layer; and the use of non-separable wavelet-transform matrices to develop truly 2D wavelet-based MRA schemes.

In its full-blown version, the envisioned MRA framework promises to significantly advance Internet theory and practice. For example, in terms of its ability to impact a more theoretical study of the Internet, it would provide a framework for unifying various Internet congestion control modeling and analysis approaches found in the current literature. On one end of the spectrum, by concentrating on the transport layer and accounting for very fine scales in space (e.g., link-to-link, host-to-host), but considering a largely trivial temporal dynamic (e.g., infinite source models), the proposed framework incorporates the scenarios treated in recent work by Low *et al.* [35, 40, 51] on the existence, uniqueness, and stability of equilibria of heterogeneous congestion control in general networks. On the other end of the spectrum, when focusing on the same transport layer and allowing for very fine scales in time (e.g., flow-level source models), but requiring an essentially trivial spatial structure (e.g., linear networks), it also captures the setup considered in recent work by Gromoll and Williams [25, 26] who study stability and heavy traffic behavior of a general stochastic flow model of congestion control for two very specific types of networks. The challenge will be to bridge the gap between these two extremes and establish similar existence, uniqueness, and stability results for models of Internet congestion control that allow for very fine scales in time **and** space. This is closely related to the problem of generalizing Kurtz’s construction to network-wide traffic matrices by (i) accounting for the spatial aspect of Internet traffic, (ii) incorporating those mechanisms of Internet congestion control that shape the behavior of network-wide traffic at the transport layer over sufficiently large time scales, and (iii) explaining features of an overall traffic matrix in terms of application-specific traffic matrices (e.g., Web traffic only, Peer-to-Peer traffic only).

From a more practical perspective, the envisioned MRA technology can also be expected to aid the development of novel and powerful tools for root-cause analyses of network failures or detection of different types of unwanted traffic (e.g., spam, botnets, worms, viruses). The ability to examine network traffic measurements in a systematic manner across many different time scales, over a variety of different spatial scales (e.g., IP address, prefix, autonomous domains), and at the different layers in the TCP/IP protocol stack suggests a holistic approach to exploiting Internet-related measurements that has been largely absent to date. In particular, it argues for tools and techniques with “drill-down” or “zoom-in” capabilities that are informed by coarse-scale representations of the data and are guided by a detailed understanding of the correlations that might exist at the different scales in time, over space, and across layers. While multi-scale approaches to, for example, network intrusion detection have been popular in the recent past, the main challenge here will be to fully exploit the multi-dimensional aspect of scale and not treat it one dimension at a time.

Acknowledgments

Matthew Roughan's participation in this work was supported by the Australian Research Council Grant DP0665427.

References

- [1] ABILENE NETWORK, <http://www.internet2.edu/abilene>.
- [2] ABRY, P. and VEITCH, D. (1998). Wavelet analysis of long-range dependent traffic. *IEEE Transactions on Information Theory*, **44**(1), pp. 2–15.
- [3] ABRY, P., TAQQU, M. S., FLANDRIN, P. and VEITCH, D. (2000). Wavelets for the analysis, estimation, and synthesis of scaling data. *Self-Similar Network Traffic and Performance Evaluation*, K. Park and W. Willinger (Editors), pp. 39–88, Wiley, New York.
- [4] ABRY, P., FLANDRIN, P., TAQQU, M. S. and VEITCH, D. (2003). Self-similarity and long-range dependence through the wavelet lens. *Long-range Dependence: Theory and Applications*, P. Doukhan, G. Oppenheim, and M. S. Taqu (Editors), pp. 527–556, Birkhäuser, Boston.
- [5] ADLER, R. J., FELDMAN, R. E., and TAQQU, M. S. (1998). *A Practical Guide to Heavy Tails: Statistical Techniques and Applications*. Birkhäuser, Boston.
- [6] ALDERSON, D., CHANG, H., ROUGHAN, M., UHLIG, S., and WILLINGER, W. (2006). The many facets of Internet topology and traffic. *Networks and Heterogeneous Media*, textbf1(4), pp. 569–600.
- [7] APPENZELLER, G., KESLASSY, I., and MCKEOWN, N. (2004). Sizing router buffers. *Computer Communication Review (Proc. of ACM/Sigcomm'04, Portland, OR)*, **34**(4), pp. 281–292.
- [8] BERAN, J. (1994). *Statistics for Long-Memory Processes*. Chapman & Hall, New York.
- [9] CAO, J., CLEVELAND, W. S., and SUN, D. X. (2002). Internet traffic tends toward Poisson and independent as the load increases. In: *Nonlinear Estimation and Classification*, C. Holmes, D. Dennison, M. Hansen, B. Yu, and B. Mallick (Editors), pp. 83–109, Springer-Verlag, New York.
- [10] COX, D. R. (1984). Long-range dependence: A review. *Statistics: An Appraisal*, H. A. David and H. T. David (Editors), pp. 55–74, Iowa State University Press, Ames, Iowa.
- [11] CROVELLA, M. E. and BESTAVROS, A. (1996). Self-similarity in World Wide Web traffic—evidence and possible causes. *Proc. ACM/Sigmetrics'96*, Philadelphia, PA, pp. 160–169.
- [12] CROVELLA, M. E. and BESTAVROS, A. (1997). Self-similarity in World Wide Web traffic—evidence and possible causes. *IEEE/ACM Transactions on Networking*, **5**, pp. 835–846.
- [13] CROVELLA, M. E. and KOLACZYK, E. (2003). Graph wavelets for spatial traffic analysis. *Proc. IEEE Infocom*.
- [14] CROVELLA, M. E., TAQQU, M. S., and BESTAVROS, A. (1998). Heavy-tailed probability distributions in the World Wide Web. *A Practical Guide to Heavy Tails: Statistical Techniques and Applications*, R. Adler, R. Feldman, and M. S. Taqu (Editors), pp. 27–53, Birkhäuser, Boston.
- [15] CROVELLA, M. E. and KRISHNAMURTHY, B. (2006). *Internet Measurements: Infrastructure, Traffic, and Applications*. J. Wiley & Sons, New York.
- [16] DAUBECHIES, I. (1992). *Ten Lectures on Wavelets*. CBMS-NSF Regional Conference Series in Applied Mathematics, Vol. 61, SIAM.

- [17] JIANG, H. and DOVROLIS, C. (2005). Why is Internet traffic bursty in short (sub-RTT) time scales? *Proc. ACM/Sigmetrics'05*, Banff, Canada.
- [18] ERRAMILI, A., NARAYAN, O., and WILLINGER, W. (1996). Experimental queueing analysis with long-range dependent packet traffic. *IEEE/ACM Transactions on Networking*, **4**(2), pp. 209-223.
- [19] ERRAMILI, A., ROUGHAN, M., VEITCH, D., and WILLINGER, W. (2002). Self-similar traffic and network dynamics. *Proceedings of the IEEE*, **90**(5), pp. 800-819.
- [20] FELDMANN, A. (2000). Characteristics of TCP connection arrivals. *Self-Similar Network Traffic and Performance Evaluation*, K. Park and W. Willinger (Editors), pp. 367-399, Wiley, New York.
- [21] FELDMANN, A., GILBERT, A. C., WILLINGER, W., and KURTZ, T. G. (1998). The changing nature of network traffic: Scaling phenomena. *Computer Communication Review*, **28**, pp. 5-29.
- [22] FELDMANN, A., GILBERT, A. C., HUANG, P., and WILLINGER, W. (1999). Dynamics of IP traffic: A study of the role of variability and the impact of control. *Proc. ACM/Sigcomm'99*, Cambridge, MA, pp. 301-313.
- [23] FLOYD, S. and PAXSON, V. (2001). Difficulties in simulating the Internet. *IEEE/ACM Transactions on Networking*, **9**:4, pp. 392-403.
- [24] FOWLER, H. J. and LELAND, W. E. (1991). Local area network traffic characteristics, with implications for broadband network congestion management. *IEEE Journal on Selected Areas in Communication*, **9**, pp. 1139-1149.
- [25] GROMOLL, H. C. and WILLIAMS, R. J. (2006). Fluid limit of a network with fair bandwidth sharing and general document size distributions. *Preprint*.
- [26] GROMOLL, H. C. and WILLIAMS, R. J. (2006). Fluid model for a data network with α -fair bandwidth sharing and general document size distributions: Two examples of stability. *This volume*.
- [27] JOO, Y., RIBEIRO, V., FELDMANN, A., GILBERT, A. C., and WILLINGER, W. (2001). TCP/IP traffic dynamics and network performance: A lesson in workload modeling, flow control, and trace-driven simulations. *Computer Communication Review*, **31**(2), pp. 25-37.
- [28] KAJ, I. and TAQQU, M. S. (2005) Convergence to fractional Brownian motion and to the Telecom process: The integral representation approach. *Preprint*.
- [29] KANNAN, J., JUNG, J., PAXSON, V., and KOKSAL, C. E. (2006). Semi-automated discovery of application session structure. *ACM/Sigcomm Internet Measurement Conference IMC'06*, Rio de Janeiro, Brazil (to appear).
- [30] KARAGIANNIS, T., MOLLE, M. and FALOUTSOS, M.. Long-range dependence: Ten years of Internet traffic modeling. *IEEE Internet Computing*, **8**(5), pp. 57-64.
- [31] KURTZ, T. G. (1996). Limit theorems for workload input models. *Stochastic Networks: Theory and Applications*, F.P. Kelly, S. Zachary and I. Ziedins (Editors), pp. 119-139, Oxford University Press, Oxford, UK.
- [32] LELAND, W. E., TAQQU, M. S., WILLINGER, W. and WILSON. D. V. (1993). On the self-similar nature of Ethernet traffic. *Proc. of ACM/Sigcomm'93*, San Francisco, CA, pp. 183-193.
- [33] LELAND, W. E., TAQQU, M. S., WILLINGER, W. and WILSON. D. V. (1994). On the self-similar nature of Ethernet traffic (Extended Version). *IEEE/ACM Transactions on Networking*, **2**, pp. 1-15.
- [34] LEVY, J. B. and TAQQU, M. S. (2000). Renewal reward processes with heavy-tailed interrenewal times and heavy-tailed rewards. *Bernoulli*, **6**(1), pp. 23-44.
- [35] LOW, S. H., PAGANINI, F., and DOYLE, J. C. (2002). Internet congestion

- control. *IEEE Control Systems Magazine*, Feb. 2002, pp. 28–43.
- [36] MANDELBROT, B. B. (1963). New methods in statistical economics. *Journal of Political Economics*, **71**, pp. 421–440.
- [37] MANDELBROT, B. B. (1969). Long-run linearity, locally Gaussian processes, H-spectra and infinite variances. *International Economic Review*, **10**, pp. 82–113.
- [38] MEDINA, A., FRALEIGH, C., TAFT, N., BHATTACHARYYA, S., and DIOT, C. (2002). A taxonomy of IP traffic matrices. *Proc. SPIE ITCOM 2002*, Boston, MA.
- [39] MIKOSCH, T., RESNICK, S., ROOTZEN, H., and STEGEMAN, A. (2002). Is network traffic approximated by stable Levy motion or fractional Brownian motion? *Annals of Applied Probability*, **12(1)**, pp. 23–68.
- [40] PAGANINI, F., WANG, Z., DOYLE, J. C., and LOW, S. H. ((2005). Congestion control for high performance, stability, and fairness in general networks. *IEEE/ACM Transactions on Networking*, **13(1)**, pp. 43–56.
- [41] PARK, K. and WILLINGER, W. (2000). *Self-Similar Network Traffic and Performance Evaluation*. J. Wiley & Sons, New York.
- [42] PAXSON, V. and FLOYD, S. (1994). Wide-area traffic: The failure of Poisson modeling. *Computer Communication Review (Proc. of ACM/Sigcomm'94, London, UK)*, **24(4)**, pp. 257–268.
- [43] PAXSON, V. and FLOYD, S. (1995). Wide area traffic: The failure of Poisson modeling. *IEEE/ACM Transactions on Networking*, **3**, pp. 226–244.
- [44] PIPIRAS, V., TAQQU, M. S., and LEVY, J. B. (2004). Slow, fast, and arbitrary growth conditions for renewal reward processes when the renewals and the rewards are heavy-tailed. *Bernoulli*, **10**, pp. 121–163.
- [45] RESNICK, S. I. (1997). Heavy tail modeling and teletraffic data. *The Annals of Statistics*, **25**, pp. 1805–1869.
- [46] ROUGHAN, M., KALMANEK, C. R. (2003). Pragmatic modeling of broadband access traffic. *Computer Communications*, **26(8)**, pp. 804–816.
- [47] ROUGHAN, M. (2005). Simplifying the synthesis of Internet traffic matrices. *Computer Communication Review*, **35**, pp. 93–96.
- [48] ROUGHAN, M., GREENBERG, A., KALMANEK, C., RUMSEWICZ, M., YATES, J., and ZHANG, Y. (2003). Experience in measuring Internet backbone traffic variability: Models, metrics, measurements, and meaning. *Proc. ITC 18*, Berlin, Germany, pp. 379–388.
- [49] SAMORODNITSKY, G. and TAQQU, M. S. (1994). *Stable Non-Gaussian Random Processes: Stochastic Models with Infinite Variance*. Chapman & Hall, London.
- [50] SOMMERS, J., BARFORD, P., DUFFIELD, N., and RON, A. (2007). Efficient network-wide SLA compliance monitoring. *Computer Communication Review (Proc. of ACM/Sigcomm'07, Kyoto, Japan)*, **39(4)** (to appear).
- [51] TANG, A., WANG, J., LOW, S. H., and CHIANG, M. (2005). Equilibrium of heterogeneous congestion control: Existence and uniqueness. *Proc. IEEE Infocom 2005*.
- [52] TAQQU, M. S., WILLINGER, W., and SHERMAN, R. (1997). Proof of a fundamental result in self-similar traffic modeling. *Computer Communication Review*, **27**, pp. 5–23.
- [53] VETTERLI, M. and KOVACEVIC, J. (1995). *Wavelets and Subband Coding*. Prentice Hall, Englewood Cliffs, NJ.
- [54] WILLINGER, W., TAQQU, M. S., and ERRAMILI, A. (1996). A bibliographical guide to self-similar traffic and performance modeling for modern high-speed

- networks. *Stochastic Networks: Theory and Applications*, F. P. Kelly, S. Zachary, and I. Ziedins (Editors), pp. 339–366, Oxford University Press, Oxford, UK.
- [55] WILLINGER, W., TAQQU, M. S., LELAND, W. E., and WILSON, D. V. (1995). Self-similarity in high-speed packet traffic: Analysis and modeling of Ethernet traffic measurements. *Statistical Science*, **10**(1), pp. 67-85.
- [56] WILLINGER, W., TAQQU, M. S., SHERMAN, R., and WILSON, D. V. (1997). Self-similarity through high-variability: Statistical analysis of Ethernet LAN traffic at the source level. *IEEE/ACM Transactions in Networking*, **5**(1), pp. 71–86.
- [57] WILLINGER, W., PAXSON, V., and TAQQU, M. S. (1998). Self-similarity and heavy tails: Structural modeling of network traffic. *A Practical Guide to Heavy Tails: Statistical Techniques and Applications*, R. Adler, R. Feldman, and M. S. Taqqu (Editors), pp. 27–53, Birkhäuser, Boston.
- [58] WILLINGER, W., ALDERSON, D., and LI, L. (2004). A pragmatic approach to dealing with high-variability in network measurements. *Proc. 2004 ACM/Sigcomm Internet Measurement Conference (IMC'04)*, pp. 88-100.
- [59] ZHANG, Y., ROUGHAN, M., LUND, C., and DONOHO, D. (2005). Estimating point-to-point and point-to-multipoint traffic matrices: An information-theoretic approach. *IEEE/ACM Transactions on Networking*, **13**, pp. 947–960.