

TECHNICAL REPORT

Multiscale Modeling and Queuing Analysis of Long-Range-Dependent Network Traffic

*Vinay J. Ribeiro, Rudolf H. Riedi, Matthew S. Crouse,
and Richard G. Baraniuk**

Department of Electrical and Computer Engineering
Rice University
6100 South Main Street
Houston, TX 77005, USA

Abstract

We develop a simple multiscale model for the analysis and synthesis of nonGaussian, long-range-dependent (LRD) network traffic loads. The wavelet transform effectively decorrelates LRD signals and hence is well-suited to model such data. However, traditional wavelet-based models are Gaussian in nature which one may at the most hope to match second order statistics of inherently nonGaussian traffic loads. Using a multiplicative superstructure atop the Haar wavelet tree, we retain the decorrelating properties of wavelets while simultaneously capturing the positivity and “spikiness” of nonGaussian traffic. This leads to a swift $O(N)$ algorithm for fitting and synthesizing N -point data sets. The resulting model belongs to the class of multifractal cascades, a set of processes with rich scaling properties which are better suited than LRD to capture burstiness. We elucidate our model’s ability to capture the covariance structure of real data and then fit it to real traffic traces. We derive approximate analytical queuing formulas for our model, also applicable to other multiscale models, by exploiting its multiscale construction scheme. Queuing experiments demonstrate the accuracy of the model for matching real data and the precision of our theoretical queuing results, thus revealing the potential use of the model for numerous networking applications. Our results indicate that a Gaussian assumption can lead to over-optimistic predictions of tail queue probability even when taking LRD into account.

*This work was supported by the National Science Foundation, grant no. MIP-9457438, by DARPA/AFOSR, grant no. F49620-97-1-0513, and by Texas Instruments. Email: {vinay, riedi, mcrouse, richb}@rice.edu. URL: www.dsp.rice.edu.

1 Introduction

Traffic models play a significant rôle in the analysis and characterization of network traffic and network performance. Accurate models capture important characteristics of traffic and enhance our understanding of these complicated signals and systems by allowing us to study the effect of various model parameters on network performance through both analysis and simulation.

One key property of modern network traffic is the presence of *long-range dependence* (LRD) which was demonstrated convincingly in the landmark paper of Leland et. al. [1]. There, measurements of traffic load on an Ethernet were attributed to *fractal* behavior or *self-similarity*, i.e., to the fact that the data “looked statistically similar” (highly variable) on all time-scales. These features are inadequately described by classical traffic models such as Markov or Poisson models. In particular, the LRD of data traffic can lead to higher packet losses than that predicted by classical queuing analysis [1, 2].

These findings were immediately followed by the development of new fractal traffic models [3–5]. *Fractional Brownian motion* (fBm), the most broadly applied fractal model, is the unique Gaussian process with stationary increments and the scaling property:

$$B(at) \stackrel{fd}{=} a^H B(t), \quad (1)$$

for all $a > 0$ with equality in the sense of finite-dimensional distributions. The parameter H , $0 < H < 1$, is known as the *Hurst parameter*. The discrete increment process $G(k) := B((k+1)\Delta) - B(k\Delta)$, called *fractional Gaussian noise* (fGn), has an autocorrelation of the form

$$r_G[k] = \frac{\sigma^2}{2} |\Delta|^{2H} (|k+1|^{2H} - 2|k|^{2H} + |k-1|^{2H}), \quad (2)$$

where Δ is a constant. Gaussianity and the strong scaling (1) enable rigorous analytical studies of queuing behavior [6–10], thus increasing the popularity of the fBm/fGn models.

Though fGn is an appropriate traffic model in some cases [11, 12], it cannot model real-

world traces with a correlation structure which is poorly matched by the rigid, restrictive (2). Indeed, convincing evidence has been produced establishing the importance of short-term correlations for buffering [13–15] and so-called relevant time scales have been discovered [14, 16, 17].

Generalizations of fBm/fGn with a more flexible correlation structure than (2) can be synthesized almost effortlessly using the amazing decorrelating capability of the *wavelet transform* [18–21]. There, independent Gaussian wavelet coefficients with variance decaying appropriately with scale form the building blocks for modeling both the long and short-term correlations of a target data set. Efficient $O(N)$ algorithms based on the tree structure of wavelet coefficients are available to synthesize N -point data sets [22, 23]. We will collect such models under the term *wavelet-domain independent Gaussian* (WIG) models.

As a consequence of their Gaussian nature, the fBm/fGn/WIG models can produce unrealistic synthetic traffic traces in certain situations. In many networking applications, for instance, we are nowhere near the Gaussian limit, in particular on small time scales. Indeed, various authors have observed heavy-tailed marginals in traffic [24, p. 364], [25]. More practically speaking, when the standard deviation of the data approaches or exceeds the mean, considerable parts of the fBm/fGn/WIG output are negative (see Figure 6(a) and (b)).

This paper has two main contributions. The first is a model for network traffic, the *multifractal wavelet model* (MWM), based on a multiplicative cascade in the wavelet domain that by design guarantees a positive output [26]. In its simplest form, the MWM is closely related to the wavelet-based construction of fBm/fGn, having the same short list of parameters (mean, variance, H). However, the MWM framework boasts the flexibility, if desired, to additionally match the short-term correlations like the WIG model. Each sample of the MWM process is obtained as a product of several positive independent random variables. The result is a positive nonGaussian, LRD process with multifractal properties. The MWM

is, thus, a more natural fit for positive arrival processes than WIG. Indeed, multifractal properties of TCP traffic have been observed first in [27] and have since shown great potential for advancing the understanding of modern networks [28, 29].

Fitting the MWM to real traffic traces results in an excellent match, far better than the WIG model, visually (see Figure 6(c)) and, as we will see, in the burstiness as measured by the multifractal spectrum, in the marginals and in the queueing behavior [26, 30]. It thus appears that the multiplicative MWM approach is more appropriate than an additive Gaussian one.

The second main contribution of this paper is a novel multiscale queueing analysis of the MWM, applicable to multiscale models in general including tree-based models like the WIG. By restricting our analysis to data at time scales of powers of two, we exploit the inherent binary tree structure of the MWM in deriving an easy-to-use and — as numerous experiments verify — close approximation to the tail queue probability, valid for any given buffer size. As a consequence, the MWM becomes viable for applications like call admission control.

After introducing wavelets and explaining the WIG model in Section 2, we describe the MWM and demonstrate the importance of the nonGaussian nature of traffic on queueing in Section 3. We then introduce the novel *multiscale queueing analysis*, which we apply to the WIG and MWM. We provide empirical evidence for the accuracy of our theoretical queueing formulas in Section 4. A tutorial introduction to multifractal cascades is found in Section 5. The proof of an instrumental lemma appears in the Appendix.

2 Classical Wavelet Models for LRD Processes

2.1 Long-range dependence

Consider a discrete-time, wide-sense stationary random process $\{X_t, t \in \mathbb{Z}\}$ with autocovariance function $r_X[k] = \text{cov}(X_t, X_{t+k})$. A change in time scale can be represented by

forming the aggregate process $X_t^{(m)}$, which is obtained by averaging X_t over non-overlapping blocks of length m and replacing each block by its mean

$$X_t^{(m)} = \frac{X_{tm-m+1} + \cdots + X_{tm}}{m}. \quad (3)$$

Denote the auto-covariance of $X_t^{(m)}$ by $r_X^{(m)}[k]$. The process X is said to exhibit LRD if its auto-covariance decays slowly enough to render $\sum_{k=-\infty}^{\infty} r_X[k]$ infinite [31]. Equivalently, $m r_X^{(m)}[0] \rightarrow \infty$ as $m \rightarrow \infty$, and the power spectrum $S_X(f)$ is singular near $f = 0$.

An important class of LRD processes are the *asymptotically second-order self-similar processes*, which are defined by the property $r_X[k] \simeq k^{2H-2}$ for some $H \in (1/2, 1)$ or, equivalently [31],

$$r_X^{(m)}[k] \rightarrow \frac{r_X^{(m)}[0]}{2} (|k+1|^{2H} - 2|k|^{2H} + |k-1|^{2H}) \quad (4)$$

as $m \rightarrow \infty$. In words, these processes “look similar” on all scales, at least from the point-of-view of second-order statistics. An example of such a process is the fGn, where the *Hurst parameter* H in (1) is exactly the scaling parameter in (4).

To estimate H by the *variance-time plot* method, we fit a straight line through the plot of an estimate of $\log \text{var}(X^{(m)})$ against $\log(m)$. More reliable estimators of H have been devised [32], in particular an unbiased one based on wavelets (see [33] and also [34]).

2.2 Wavelet transform

The discrete wavelet transform provides a multiscale signal representation of a one-dimensional random signal $C(t)$ in terms of shifted and dilated versions of a prototype bandpass wavelet function $\psi(t)$ and shifted versions of a lowpass scaling function $\phi(t)$ [35, 36]. For special choices of the wavelet and scaling functions, the atoms

$$\psi_{j,k}(t) := 2^{j/2} \psi(2^j t - k), \quad \phi_{j,k}(t) := 2^{j/2} \phi(2^j t - k), \quad j, k \in \mathbb{Z} \quad (5)$$

form an orthonormal basis, and we have the signal representation [36]

$$C(t) = \sum_k U_{J_0,k} \phi_{J_0,k}(t) + \sum_{j=J_0}^{\infty} \sum_k W_{j,k} \psi_{j,k}(t). \quad (6)$$

Here the wavelet coefficients $W_{j,k}$ and the scaling coefficients $U_{j,k}$ are given by

$$W_{j,k} := \int C(t) \psi_{j,k}(t) dt, \quad U_{j,k} := \int C(t) \phi_{j,k}(t) dt. \quad (7)$$

Without loss of generality, we will assume $J_0 = 0$.

In this representation, k indexes the spatial location of analysis and j indexes the *scale* or resolution of the wavelet analysis — larger j corresponds to higher resolution and $j = 0$ indicates the coarsest scale or lowest resolution of analysis. In practice, we work with a sampled or finite-resolution representation of $C(t)$, replacing the semi-infinite sum in (6) with a sum over a finite number of scales $0 \leq j \leq n - 1$, $n \in \mathbb{Z}_+$.

In this paper, we restrict our attention to the simplest wavelet system, the *Haar*. The Haar scaling and wavelet functions are given by (see Figure 2(a))

$$\phi(t) = \begin{cases} 1, & 0 \leq t < 1 \\ 0, & \text{else} \end{cases} \quad \text{and} \quad \psi(t) = \begin{cases} 1, & 0 \leq t < 1/2 \\ -1, & 1/2 \leq t < 1 \\ 0, & \text{else.} \end{cases} \quad (8)$$

Since $\phi_{j,k}(t)$ is a rectangular function, the Haar scaling coefficients $U_{j,k}$ (7) represent the *local mean values* of the signal in the time intervals $[k2^{-j}, (k+1)2^{-j}]$ and thus form a discrete-time approximation of $c(t)$ at resolution j . By design the support of $\phi_{j,k}(t)$ are nested within each other. This makes it natural to use a binary tree (Figure 2(b)) to display the relationship between coefficients $U_{j,k}$. Nodes at lower horizontal levels in the tree correspond to representations of the signal at finer resolutions.

The Haar wavelet transform of a signal can be computed recursively starting from its finest-scale scaling coefficients via [36]

$$U_{j-1,k} = 2^{-1/2}(U_{j,2k} + U_{j,2k+1}), \quad W_{j-1,k} = 2^{-1/2}(U_{j,2k} - U_{j,2k+1}). \quad (9)$$

This corresponds to moving up the binary tree and storing in the Haar wavelet coefficients $w_{j,k}$ the detail information lost while going from fine to coarse resolutions (Figure 2(b)).

The inverse Haar wavelet transform is computed via

$$U_{j,2k} = 2^{-1/2}(U_{j-1,k} + U_{j-1,k}) \quad \text{and} \quad U_{j,2k+1} = 2^{-1/2}(U_{j-1,k} - W_{j-1,k}) \quad (10)$$

and is equivalent to moving down the scaling coefficient tree to get finer representations of the signal. It is easily shown that the forward and inverse Haar wavelet transforms of an N -point signal can be computed in $O(N)$ operations, using (9) and (10) respectively.

We introduce three different processes: the continuous-time signal $C(t)$, its integral $D(t)$, and a discrete-time approximation $C^{(n)}[k]$ to $C(t)$. These three signals are related by

$$C^{(n)}[k] := \int_{k2^{-n}}^{(k+1)2^{-n}} C(t) dt = D\left((k+1)2^{-n}\right) - D\left(k2^{-n}\right). \quad (11)$$

In this paper, $C^{(n)}[k]$ and $D(t)$ will play rôles analogous to fGn and fBm, respectively.

For notational simplicity, we will assume that both $C(t)$ and $D(t)$ live on $[0, 1]$ and that $C^{(n)}[k]$ is a length- 2^n discrete-time signal. Thus, there is only one scaling coefficient $U_{0,0}$ in (6), that is, a single tree of scaling coefficients. (A more general case with multiple scaling coefficients at the coarsest scale is treated in [37].) $C^{(n)}[k]$ relates directly to the finest-scale scaling coefficients:

$$C^{(n)}[k] = 2^{-n/2}U_{n,k}, \quad k = 0, 1, \dots, 2^n - 1. \quad (12)$$

We will focus on modeling $C^{(n)}[k]$ in this paper.

2.3 Wavelet-domain Independent Gaussian (WIG) model

Wavelets serve as an approximate Karhunen-Loève or decorrelating transform for fBm [18], fGn, and more general LRD signals [23]. Hence, the difficult task of modeling these highly correlated signals in the time domain reduces to a simple one of modeling them approximately by an uncorrelated process in the wavelet domain.

The WIG model synthesizes a Gaussian LRD process by generating the parent node $U_{0,0}$ of the scaling coefficient tree as a Gaussian random variable with mean equal to the sample mean of the process, and by generating the wavelet coefficients as independent, i.e. uncorrelated, zero-mean Gaussian random variables which are identically distributed within scale according to $W_{j,k} \sim N(0, \sigma_j^2)$, with σ_j^2 the wavelet-coefficient variance at scale j [18–22].

Scaling coefficients at finer scales on the tree are then recursively computed through (10) until the finest scale scaling coefficients $U_{n,k}$ and hence the required signal $C_{\text{WIG}}^{(n)}[k]$ are obtained. The result is a fast $O(N)$ algorithm for generating a length- N signal (see Figure 2(c)).

An attractive feature of the WIG model is its flexibility in matching different correlation structures of LRD processes. A power-law decay for the σ_j^2 's leads to approximate wavelet synthesis of fBm or fGn [18, 20]. However, while network traffic may exhibit correlations consistent with fBm or fGn asymptotically at large scale, it may have short-term correlations that vary considerably from pure fBm or fGn scaling. Such LRD processes can be modeled by setting σ_j^2 to match the measured or theoretical variances of the wavelet coefficients of the desired process [22]. Thus, for a length- N signal, the WIG is characterized by approximately $\log_2(N)$ (the number of time scales) parameters.

The WIG is an *additive model*, because we can express the signal $C_{\text{WIG}}^{(n)}[k]$ directly as a sum of independent random variables. First, we need some notation. Each shift k at scale n has a unique binary representation $k = \sum_{i=0}^{n-1} k'_i 2^{n-1-i}$ where each $k'_i \in \{0, 1\}$. Letting $k_n = k$ and $k_{i-1} = k_i \text{ div } 2$ we have $k_i = 2 \cdot k_{i-1} + k'_{i-1} = \sum_{j=0}^{i-1} k'_j 2^{i-1-j}$. The shifts k_i correspond to the ancestors of k at scale i and so we can write

$$C_{\text{WIG}}^{(n)}[k] = 2^{-n/2} U_{n,k} = 2^{-n} \left(U_{0,0} + \sum_{i=0}^{n-1} (-1)^{k'_i} 2^{i/2} W_{i,k_i} \right), \quad (13)$$

This result can be derived by iteratively applying (10).

The WIG model is Gaussian by construction, but network traffic signals (such as loads and interarrival times) can be highly nonGaussian (Figure 6). Not only are these signals strictly non-negative, but also they exhibit “spiky” behavior corresponding to marginals with tails that are somewhat heavier than the Gaussian. Moreover, this spikyness is poorly captured by LRD. We seek a more accurate marginal characterization for these spiky, non-negative LRD processes, yet wish to retain the decorrelating properties of wavelets and the simplicity of the WIG model. This leads us to introduce our multiplicative model and the concept of multifractals in the following sections.

3 Multifractal Wavelet Model

3.1 Haar wavelet transform and positive signals

In order to model non-negative signals using the Haar wavelet transform, we must constrain the scaling and wavelet coefficient values to ensure that $C(t)$ in (6) is non-negative. While cumbersome for a general wavelet system,¹ these conditions are simple for the Haar system.

Since the Haar scaling coefficients $u_{j,k}$ represent the local mean of the signal at different scales and shifts, they are non-negative if and only if the signal itself is non-negative, that is, $C(t) \geq 0 \Leftrightarrow u_{j,k} \geq 0, \forall j, k$. This condition leads us directly to a set of constraints on the Haar wavelet coefficients. Combining (10) with the constraint $u_{j,k} \geq 0$, we obtain the condition

$$C(t) \geq 0 \Leftrightarrow |w_{j,k}| \leq u_{j,k}, \quad \forall j, k. \quad (14)$$

3.2 MWM model

The positivity constraints (14) inspire a very simple multiscale, multiplicative signal model for positive processes. In the *multifractal wavelet model* (MWM) [37] we compute the wavelet coefficients recursively by

$$W_{j,k} = A_{j,k} U_{j,k}, \quad (15)$$

where the $A_{j,k}$'s are independent random variables supported on the interval $[-1, 1]$.

The MWM synthesizes a data trace in a manner similar to the WIG. After generating the coarsest scale scaling coefficient $U_{0,0}$ and the multipliers $A_{j,k}$, the MWM generates scaling coefficients at finer scales of the scaling coefficient tree recursively using (10) and (15), that is (see Figure 3)

$$U_{j,2k} = 2^{-1/2}(1 + A_{j+1,k})U_{j-1,k}, \quad U_{j,2k+1} = 2^{-1/2}(1 - A_{j+1,k})U_{j-1,k}, \quad (16)$$

until the finest scale has been reached. The MWM is a *multiplicative model* because we can express the signal $C_{\text{MWM}}^{(n)}[k]$ directly as a product (or cascade) of independent random

¹The conditions are straightforward also for certain biorthogonal wavelet systems.

multipliers $1 \pm A_{j,k}$. Using the notation introduced in Section 2.3, we have

$$C_{\text{MWM}}^{(n)}[k] = 2^{-n/2} U_{n,k} = 2^{-n} U_{0,0} \prod_{i=0}^{n-1} (1 + (-1)^{k_i} A_{i,k_i}), \quad (17)$$

which should be compared with (13).²

It is easily shown that the total cost for computing N MWM signal samples is $O(N)$. In fact, synthesis of a trace of length 2^{18} data points takes just seconds of workstation cpu time. See [38] for a similar model to the MWM used as an intensity prior for wavelet-based image estimation.

3.3 Model parameters

We choose the multipliers $A_{j,k}$ to be symmetric about 0 and identically distributed within scale; it is easily shown that these two conditions are necessary for the $C_{\text{MWM}}^{(n)}$ process to be first-order stationary [37]. Due to its flexible shape (see Figure 3(b)), compact support and tractability to closed-form calculations, we choose the *symmetric beta distribution*³ [39], $\beta_{-1,1}(p_j, p_j)$ (see Figure 3(b)) for the $A_{j,k}$'s, with p_j the beta parameter at scale j .

A precise model of $U_{0,0}$ requires a strictly non-negative probability density function to ensure the non-negativity of the MWM output. In our simulation experiments we choose $U_{0,0} \sim \beta_{0,M}(p_{-1}, p_{-1})$ with $M \geq 0$.

3.4 Model training

Since the variance of a random variable $A \sim \beta_{-1,1}(p, p)$ is

$$\text{var}[A] = 1/(2p + 1), \quad (18)$$

we obtain from (15)

$$\frac{\text{var}(W_{j-1,k})}{\text{var}(W_{j,k})} = \frac{2 \text{var}[A_{j-1,k}]}{\text{var}[A_{j,k}] (1 + \text{var}[A_{j-1,k}])} = \frac{2p_j + 1}{p_{j-1} + 1}. \quad (19)$$

²It is tempting to exponentiate (13) in order to obtain a process with a multiplicative structure as (17). However, it should be noted that such an approach bares various difficulties. First, the exponential of a zero mean variable ($W_{j,k}$) is not of mean one ($1 \pm A_{j,k}$). Second, it is hard to control correlation structure under exponentiation.

³We denote a beta random variable with support $[a, b]$ by $\beta_{a,b}$

Table 1: Comparison of the tree-based WIG and MWM models. For approximating a signal with a strict fGn covariance structure, both the WIG and MWM require only three parameters (mean, variance, and H).

WIG	MWM
Additive	Multiplicative
Gaussian	Asymptotically Lognormal
LRD matched	LRD matched
Independent wavelet coefficients	Independent multipliers
Monofractal	Multifractal
$\log_2 N + 2$ parameters	$\log_2 N + 2$ parameters
$O(N)$ synthesis	$O(N)$ synthesis

Thus, to model a target process with the MWM, we select the p_j 's to match the signal's theoretical wavelet-domain energy decay. Or, given training data, we can select the parameters to match the sample variances of the wavelet coefficients as a function of scale. With one beta parameter per wavelet scale, the MWM has approximately $\log_2 N$ parameters for a trace of length N . Distributions with more parameters (e.g., discrete distributions or mixtures of betas) could be used to capture high-order data moments at a cost of increased model complexity [37]. See Table 1 for a comparison of the WIG and MWM properties.

To complete the modeling, we must choose the parameters p_0 , p_{-1} and M of the model. From (15) and (18) we obtain $(2p_0 + 1)\text{var}(W_{0,0}) = \mathbb{E}[U_{0,0}^2]$, which allows us to calculate p_0 from estimates of $\mathbb{E}[U_{0,0}^2]$ and $\text{var}(W_{0,0})$. The parameters p_{-1} and M of $U_{0,0}$ are chosen using estimates of $\mathbb{E}(U_{0,0})$ and $\text{var}(U_{0,0})$.

3.5 Experimental results

The capability of both models, the WIG and the MWM, were tested using two real data traces which are well studied in the literature. The first trace (LBL-TCP-3) contains two hours wide-area TCP traffic between the Lawrence Berkeley Laboratory and the rest of the world in 1994 [40] and the second trace (BC-pAug89) is one of the celebrated Ethernet data traces collected at Bellcore Morristown Research and Engineering facility in 1989 [1]. The LBL-TCP-3 trace that we use is of bytes arriving in time intervals of 6ms and the BC-pAug89

trace is of bytes per 2.6ms. To model the data, we use estimates of the σ_j^2 at the 15 finest dyadic scales where there is sufficient data to obtain good estimates.

Figure 6(c) demonstrates that the MWM produces positive “spiky” data akin to the real traffic; to the contrary with the WIG model. Also the marginals of the MWM traces are a much better match to the LBL-TCP-3 trace than those of the WIG. Figure 4 displays a comparison at three different aggregation levels

As expected the correlation structure is well matched by both, the WIG and the MWM, since these models are both designed to do so.⁴ For evidence we refer to the variance-time plots of Figure 5 which were obtained by averaging the empirical variance-time plots of 32 independent realizations of the models.

To study the impact of the nonGaussian nature of the real data on queuing, we conduct queuing experiments. We consider an infinite-length single-server queue with link capacity 800 bytes/unit time. In Figure 6, observe that the MWM traces closely match the queuing behavior of the real data traces while the WIG traces do not⁵. We conclude that a Gaussian assumption for spiky nonGaussian traffic can lead to over-optimistic predictions of tail queue probabilities. To gain further insight as to why this is the case, we are motivated to perform a theoretical queuing analysis of the WIG and MWM.

4 Multiscale Queuing Analysis

Queuing analysis is fundamental to network engineering. Buffer dimensioning in routers and call admission control are but two of the many crucial areas in networking research that rely on an accurate characterization of the queuing behavior of data traffic.

The discovery of LRD in traffic has created a challenging new area of research in queuing theory. Analytical studies prove that an infinite-length buffer with constant service rate

⁴To reduce the parameters of both, the WIG and the MWM, we could fit the VTP of the real data with say a cubic polynomial. Then, we would express the variances σ_j^2 in terms of the 4 polynomial coefficients instead of matching the variance exactly on all scales.

⁵In all experiments in this paper, confidence intervals plotted correspond to a confidence level of 95%.

fed with traffic loads from fGn-based models has a tail queue distribution which decays asymptotically like a Weibullian law

$$\mathbf{P}[Q > b] \simeq \exp(-\delta b^{2-2H}). \quad (20)$$

Here, δ is a positive constant that depends on the service rate of the queue [7, 8]. Clearly, (20) reveals that the decay of the tail queue distribution for fGn with $H > 1/2$ is much slower than the exponential decay predicted by SRD classical models [2] which correspond to the case $H = 1/2$. In spite of this result, there is still an ongoing discussion on the effect of LRD on queuing, with researchers arguing both for and against its importance [14–17, 41, 42].

In this section, we present an approach to queuing analysis which is particularly adapted to *multiscale representations* of signals and processes. More precisely, exploiting the inherent binary tree structure of the Haar scaling coefficients of both traffic models, the WIG and the MWM, we derive approximate formulas for their tail queue probability. Doing so, our queuing formulas

- are applicable to tree-based / multiresolution models in general,
- they are valid for any queue size, unlike (20) which is an asymptotic result,
- they capture more complicated correlation structures than the mere asymptotic LRD exponent H ,
- and moreover, they involve the entire distributions of data at multiple time resolutions and not only the second order statistics.

Finally, we demonstrate experimentally that our theoretical results — which involve some approximations – are in good agreement with the empirical tail queue behavior of both the WIG and the MWM.

4.1 Queue size and multiple time scales

Consider a discrete time random process, the traffic load L_i , $i \in \mathbb{Z}$ which we think of as entering an infinite buffer single server queue with constant link capacity c . Let Q_i represent the queue size at time instant i . Denote by K_r the aggregate traffic arriving between time instants $-r + 1$ and 0

$$K_r = \sum_{i=-r+1}^0 L_i. \quad (21)$$

In the sequel, we refer to K_r as representing the data at time-scale r . We set $K_0 := 0$. Using Lindley's equation [43], it is easily shown that

$$Q_0 = \max(Q_{-r} + K_r - rc, K_{r-1} - (r-1)c, \dots, K_0). \quad (22)$$

Since $Q_{-r} \geq 0$ for all r , we must have $Q_0 \geq \sup_{r \in \mathbb{Z}_+} (K_r - rc)$. Denoting by $-t$ the last instant the queue was empty before time instant 0 (we set $-t = 0$ if $Q_0 = 0$), we obtain $Q_0 = K_t - tc \leq \sup_{r \in \mathbb{Z}_+} (K_r - rc)$. Thus if the queue was empty at some time in the past,

$$Q_0 = \sup_{r \in \mathbb{Z}_+} (K_r - rc). \quad (23)$$

We will study the quantity \tilde{Q}_0 which is obtained by restricting the supremum in (23) to time scales which appear naturally in a multiscale representation, i.e. the dyadic time scales:

$$\tilde{Q}_0 := \sup_{m \in \{0, \dots, n\}} (K_{2^m} - c2^m). \quad (24)$$

Clearly, $\tilde{Q}_0 \leq Q_0$ and $\mathbf{P}(Q_0 > b) \geq \mathbf{P}(\tilde{Q}_0 > b)$.

The first approximation of our analysis is:⁶ $\mathbf{P}(Q_0 > b) \approx \mathbf{P}(\tilde{Q}_0 > b)$. To justify this, we require the notion of a *critical time scale* (CTS) [14, 16, 17]. The CTS is defined as $r^* = \arg \sup_{r \in \mathbb{Z}_+} \mathbf{P}(K_r - cr > b)$ and the *critical time-scale queue* (CTSQ) as $\text{CTSQ}(b) := \mathbf{P}(K_{r^*} - cr^* > b)$. It has been shown that $\text{CTSQ}(b) \approx \mathbf{P}(Q_0 > b)$ [14, 16, 17].

Similarly, we introduce now the *critical dyadic time-scale* (CDTS) as $r_d^* = \arg \sup_{r \in \{2^m, m \in \{0, \dots, n\}\}} \mathbf{P}(K_r - cr > b)$ and the *critical dyadic time-scale queue* (CDTSQ) as $\text{CDTSQ}(b) := \mathbf{P}(K_{r_d^*} - cr_d^* > b)$. Clearly, $\text{CDTSQ}(b) \leq \mathbf{P}(\tilde{Q}_0 > b) \leq \mathbf{P}(Q_0 > b)$.

⁶Here, \approx denotes that two quantities are approximately equal.

We present the following two heuristic arguments for our approximation:

1. If $2^n > r^*$, we are justified in neglecting time scales $r > 2^n$.
2. Dyadic time scales, though a small subset of \mathbb{Z}_+ , span the entire range of time scales. This ensures that $\text{CDTSQ} \approx \text{CTSQ}$, provided $2^n > r^*$. In fact, queuing experiments with synthetic WIG traces corresponding to an fGn correlation structure demonstrate that the CTSQ, which requires the distributions of data at all time scales, is almost equal to the CDTSQ that requires the distributions only at dyadic time scales (Figure 7). Thus, the CDTSQ appears to be a good substitute for CTSQ in on-line applications.

We observe from Figure 7 that the CTSQ is, however, not very close to the empirical tail queue probability of the WIG model. We are thus motivated to find a better approximation to the tail queue probability than the CDTSQ, which we do in the next section.

4.2 Queuing analysis

4.2.1 Queuing formula for tree-based multiscale models

In this section, we develop a new multiscale approach to queuing analysis. We derive an approximate formula for the tail queue probability of tree-based multiscale models in general⁷, including the WIG and MWM.

Performing an exact queuing analysis of tree-based models like the WIG and MWM is very complicated because their binary tree naturally produces a process that is not *strictly stationary* [37]. We would thus expect the distribution of the queue size to vary with time. For an illustration notice that in Figure 2(b) the neighboring nodes $U_{j+2,4k}$ and $U_{j+2,4k+1}$ share a parent node $U_{j+1,2k}$ at scale $j+1$ while the nodes $U_{j+2,4k+1}$ and $U_{j+2,4k+2}$ do not.

The second approximation in our analysis is to equate the tail queue probability of the models at the last instant $2^n - 1$ to the empirical tail queue probability. In other words, we choose $L_i = C[2^n - 1 + i]$ ($i = -1, -2, \dots, -2^n + 1$). The Haar scaling coefficients on the

⁷The analysis can be used for models not based on trees, but with an explicit relationship between data at different time scales

branch linking $U_{0,0}$ and $U_{n,2^n-1}$, in other words, the right edge of the tree of Figure 2(b), are then related to the quantities K_{2^m} (21) by

$$K_{2^{n-i}} = 2^{-i/2} U_{i,2^i-1}, \quad \text{for } i = 0, \dots, n. \quad (25)$$

We will later demonstrate through experiments that this approximation results in a queuing formula that closely matches the empirical tail queue probability (Section 4.3).

For the ease of notation let us denote \tilde{Q}_0 and Q_0 of Section 4.1 by \tilde{Q} and Q respectively.

Let us formulate now a queuing analysis where we assume only knowledge on the *multiresolution representation* of the arriving workload L_i . Let E_i denote the event $\{K_{2^{n-i}} < b + c2^{n-i}\}$. The following Lemma simplifies our analysis, the proof of which is given in the Appendix.

Lemma: *Assume that the events E_i are of the form $E_i = \{S_i < b_i\}$, where $S_i = X_0 + \dots + X_{i-1}$ for $1 \leq i \leq n$ and where X_0, \dots, X_n are independent, otherwise arbitrary random variables. Then for $1 \leq i \leq n$*

$$\mathbf{P}(E_i | E_{i-1}, \dots, E_0) \geq \mathbf{P}(E_i).$$

Given the Lemma we have

$$\begin{aligned} \mathbf{P}(\tilde{Q} > b) &= 1 - \mathbf{P}(\tilde{Q} < b) = 1 - \mathbf{P}(\cap_{i=0}^n E_i) \quad \text{from (24)} \\ &= 1 - \mathbf{P}(E_0) \prod_{i=1}^n \mathbf{P}(E_i | E_{i-1}, \dots, E_0) \\ &\leq 1 - \prod_{i=0}^n \mathbf{P}(E_i). \end{aligned} \quad (26)$$

We thus arrive at an upper bound approximation of $\mathbf{P}(\tilde{Q} > b)$ which would be exact if the events E_i were independent. We call this approximation the *multiscale queue* (MSQ), that is,

$$\text{MSQ}(b) := 1 - \prod_{i=0}^n \mathbf{P}(E_i). \quad (27)$$

Thus, intuitively, the MSQ assumes that dyadic time scales

- capture the effect of all time scales on the queue size and that
- they are sufficiently “far” apart for the events E_i to be considered independent.

Recall from Section 4.1 that $\tilde{Q} \leq Q$. This implies that $\mathbf{P}(\tilde{Q} > b) \leq \mathbf{P}(Q > b)$, which means that MSQ is an upper bound of a lower bound on $\mathbf{P}(Q > b)$. Our queuing approximation is thus

$$\boxed{\mathbf{P}(Q > b) \approx MSQ(b) := 1 - \prod_{i=0}^n \mathbf{P}(K_{2^{n-i}} < b + c2^{n-i})}. \quad (28)$$

Note that only multiscale marginals enter in (28).

4.2.2 Queuing analysis of the WIG

For the WIG, on choosing $X_0 := U_{0,0}$ and $X_i := -2^{i/2}W_{i,2^{i-1}}$ we obtain from (25)

$$K_{2^{n-i}} = 2^{-i} \left(U_{0,0} + \sum_{j=0}^{i-1} X_j \right) = 2^{-i} S_i. \quad (29)$$

Setting $b_i = b2^i + c2^n$, we observe that the WIG satisfies the conditions of the Lemma.

Since for the WIG $K_{2^{n-i}}$ is Gaussian, the probability $\mathbf{P}(E_i)$ can be computed from the cumulative distribution of a Gaussian distribution [39].

4.2.3 Queuing analysis of the MWM

Denoting $A_{j,2^{j-1}}$ by A_j , (25) reduces to

$$K_{2^{n-i}} = U_{0,0} \prod_{j=0}^{i-1} (1 - A_j) / 2. \quad (30)$$

The event E_i is thus

$$\begin{aligned} E_i &= \{K_{2^{n-i}} < b + c2^{n-i}\} = \{U_0 \prod_{j=0}^{i-1} (1 - A_j) < b2^i + c2^n\} \\ &= \left\{ \log(U_0) + \sum_{j=0}^{i-1} \log(1 - A_j) < \log(b2^i + c2^n) \right\}. \end{aligned} \quad (31)$$

By setting $X_0 := \log(U_0)$, $X_i := \log(1 - A_i)$ and $b_i := \log(b2^i + c2^n)$ we see that the lemma applies to the MWM. Consequently, we use (28) to approximate $\mathbf{P}(Q > b)$ for the MWM.

For the MWM, obtaining $\mathbf{P}(E_i)$ is not as straightforward as for the WIG. If $U_{0,0}$ is equal to a constant M times the random variable $\beta_{0,1}(p_{-1}, q_{-1})$, then from (30) K_{2^n-i} is M times several independent $\beta_{0,1}$ random variables. We approximate K_{2^n-i}/M by a beta random variable, $\beta_{0,1}(d_i, e_i)$, using Fan's approximation [39, 44]. Thus, if $(1 - A_j)/2 \sim \beta_{0,1}(p_j, q_j)$ then

$$d_i = S(T - S^2)^{-1}(S - T) \quad \text{and} \quad e_i = (1 - S)(T - S^2)^{-1}(S - T), \quad (32)$$

where

$$S = \prod_{j=-1}^{i-1} \frac{p_j}{p_j + q_j} \quad \text{and} \quad T = \prod_{j=-1}^{i-1} \frac{p_j(p_j + 1)}{(p_j + q_j)(p_j + q_j + 1)}. \quad (33)$$

This approximation matches the mean and variance of the actual distribution of K_{2^n-i} exactly and closely approximates the first 10 moments [44]. We thus use the cumulative distribution of a β random variable to calculate $\mathbf{P}(E_i)$, for which several approximations are available [39].

4.3 Validation of the MSQ

From Figure 6 we observe that the MSQ gives a close approximation the empirical tail queue probability of the WIG and MWM and is closer than the CDTSQ. Further experiments with synthetic traces with an fGn correlation structure confirm this result (Figures 7 and 8).

Since the MSQ uses the entire distribution of data at multiple time scales and not just the variance, we conclude that matching only the variance-time plot of heavy-tailed spiky data with a Gaussian model can lead to optimistic predictions of tail queue probability.

5 MWM is a Cascade

The MWM is closely with the theory of multiplicative cascades. Cascades provide a natural framework for producing positive “bursty” processes and offer greater flexibility and richer scaling properties than fractal models such as fGn and fBm. Closely related to cascades is the powerful theory of *multifractals*, a statistical tool for measuring “burstiness” which is much more to the point than LRD which merely measures “high variability”.

5.1 Cascades

The backbone of a cascade is a construction where one starts at a coarse scale and develops details of the process on finer scales iteratively in a multiplicative fashion. The MWM, e.g., is a multiplicative cascade: as (16) and (17) reveal we may write (see Figure 9(a))

$$C_{\text{MWM}}^{(n)}[k] = 2^{-n} M_0^0 \prod_{i=1}^n M_{i, k_i}, \quad \text{with} \quad M_{k_i}^i = \frac{(1 + (-1)^{k_i-1} A_{i-1, k_{i-1}})}{2}. \quad (34)$$

This construction procedure naturally results in a process that “sits” just above the zero line and emits occasional positive jumps or spikes. In contrast, additive self-similar models such as fGn and the WIG “hover” around the mean with occasional outbursts in both positive and negative directions.

5.2 Multifractal analysis

Intuitively, multifractal analysis measures the frequency with which bursts of different strengths occur in a signal. Consider a positive process $Y(t)$. The strength of the burst of Y at time t , also called the degree of *Hölder continuity*, can be characterized by

$$\alpha(t) = \lim_{k_n 2^{-n} \rightarrow t} \alpha_{k_n}^n \quad \text{where} \quad \alpha_{k_n}^n := -\frac{1}{n} \log_2 |Y((k_n + 1)2^{-n}) - Y(k_n 2^{-n})| \quad (35)$$

where $k_n 2^{-n} \rightarrow t$ means that $t \in [k_n 2^{-n}, (k_n + 1)2^{-n})$ and $n \rightarrow \infty$. The smaller the $\alpha(t)$, the larger the increments of Y around time t , and the “burstier” it is at time t . The frequency of occurrence of a given strength α , can be measured by the *multifractal spectrum*:

$$f(\alpha) := \lim_{\varepsilon \rightarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log_2 \#\{k_n = 0, \dots, 2^n - 1 : \alpha_{k_n}^n \in (\alpha - \varepsilon, \alpha + \varepsilon)\}. \quad (36)$$

By definition, f takes values between 0 and 1 and is often shaped like a \cap and concave. The smaller the $f(\alpha)$, the “fewer” points t will exhibit $\alpha(t) \approx \alpha$. If $\bar{\alpha}$ denotes the value $\alpha(t)$ assumed by “most” points t , then $f(\bar{\alpha}) = 1$. See Figure 9 for the multifractal spectrum of the LBL-TCP-3 data set and of synthetic MWM data. We observe that the MWM captures the spectrum of the real data except for large values of α . This means that the MWM does not generate as many small values as the signal possesses.

5.3 Multifractal spectrum and higher-order moments

Though (36) gives us a simple measure of burstiness in data, in practice it is impossible to compute the right side of (36). However, $f(\alpha)$ can be obtained through the use of high and low-order moments of the signal $Y(t)$.

Define the *partition function* that captures the scaling of different moments of Y as

$$T(q) := \lim_{n \rightarrow \infty} \frac{1}{-n} \log_2 \mathbb{E}[S_n(q)], \quad (37)$$

with

$$S_n(q) := \sum_{k_n=0}^{2^n-1} \left| Y((k_n+1)2^{-n}) - Y(k_n 2^{-n}) \right|^q = \sum_{k_n=0}^{2^n-1} 2^{-qn\alpha_{k_n}^n}. \quad (38)$$

The multifractal spectrum $f(\alpha)$ and $T(q)$ are closely related, as the following hand-waving argument shows. Grouping in the sum $S_n(q)$ of (38) the terms behaving as $\alpha_{k_n}^n \approx \alpha$, and using (36) we get

$$S_n(q) = \sum_{\alpha} \sum_{\alpha_n \sim \alpha} (2^{-n\alpha})^q \approx \sum_{\alpha} 2^{nf(\alpha)} 2^{-nq\alpha} \approx 2^{-n \inf_{\alpha} (q\alpha - f(\alpha))}. \quad (39)$$

We conclude that we must “expect” $T(q)$ to equal $\inf_{\alpha} (q\alpha - f(\alpha))$, the so-called *Legendre transform* of $f(\alpha)$. For the special case of an MWM process, i.e., $Y = D$ (see Section 2.2 for the definition of D), it can be shown (see [45]) that the inverse relation holds, called the *multifractal formalism*

$$f(\alpha) = T^*(\alpha) := \inf_q (q\alpha - T(q)). \quad (40)$$

In order to estimate $T(q)$ from a data set, it is customary to use the approximation $2^{-nT(q)} \approx S_n(q)$. For the MWM this is equivalent to

$$2^{-jT(q)} \approx \sum_{k=0}^{2^j-1} |2^{-j/2} U_{j,k}|^q. \quad (41)$$

The slope of a linear fit of $\log S_{(j)}(q)$ against j will give $T(q)$.

For the MWM, assuming the moments of the multipliers M_{i,k_i} converge to a limiting random variable $M \sim \beta_{0,1}(p, p)$, we find

$$T_D(q) = -1 - \log_2 \mathbb{E}[M^q] = \begin{cases} -1 - \log_2 \frac{\Gamma(p+q)\Gamma(2p)}{\Gamma(2p+q)\Gamma(p)} & \text{if } q > -p \\ -\infty & \text{if } q \leq -p. \end{cases} \quad (42)$$

For the self-similar fBm,

$$T_{\text{fBm}}(q) = \begin{cases} qH - 1 & \text{for } q > -1, \\ -\infty & \text{for } q \leq -1. \end{cases} \quad (43)$$

On taking the Legendre transform of T_{fBm} we observe that fBm possesses only one degree of “burstiness” ($\alpha(t) = H$) which is omnipresent. Consequently, fBm (or its increments process fGn) cannot capture the complicated multifractal behavior or “burstiness” of real data like the LBL-TCP-3 trace (Figure 9).

6 Conclusions

The MWM provides a new multiscale tool for synthesis of nonGaussian LRD traffic. Computations involving the MWM are extremely efficient — synthesis of a trace of N sample points requires only $O(N)$ computations. In fact, synthesis of even long 2^{18} point data sets takes just seconds of workstation cpu time. The parameters of the MWM, numbering approximately $\log N$, are identical in number to the WIG model and are simple enough to be either inferred from observed data or chosen a priori. We can reduce the number of parameters further by developing a parametric characterization of the wavelet energy decay across scale. A linear fit of the variance-time plot, e.g., would lead to the analogue of what the fGn is for Gaussian LRD processes. Higher order polynomial fits would provide better matches at the cost of more parameters, until in the extreme the number of parameters equals the number of scales as we have chosen to do here.

With the MWM, we have been able to fit actual traffic traces, and have developed preliminary queueing results that demonstrate that modeling heavy-tailed spiky data with Gaussian models can lead to over-optimistic predictions of tail queue probability.

We derived a closed form approximate queueing formula for the MWM that uses the cumulative distributions of data at dyadic time scales and demonstrated its accuracy through experiments. As a consequence, the versatile MWM model is now viable for numerous networking applications including call admission control.

Future research will aim at making the MWM practicable for prediction. The parameters of the MWM could also be used to capture the effect of different protocols on shaping data flow, e.g. the protocols over TCP such as FTP and HTTP. In short, the use of the MWM in real-time network protocols and control algorithms seems very promising.

Appendix

Lemma: *Assume that X_i ($i = 0, \dots, n$) are independent, but otherwise arbitrary random variables and let $S_i = X_0 + \dots + X_i$. Let b_i be arbitrary numbers. Then*

$$\mathbf{P}(S_i < b_i | S_l < b_l, l = 0, \dots, i-1) \geq \mathbf{P}(S_i < b_i).$$

Proof

Let us first spell out some notation. By f_Z and F_Z we denote the probability density function and cumulative density function, respectively, of a random variable Z . Furthermore, we denote by $F_{Z|E}(z)$ the cumulative density function of Z knowing the event E .

For convenience, let us write $E_i := \{S_i < b_i\}$ for short, and let us introduce the auxiliary random variables $Y_0 := Z_0 := S_0 = X_0$ and

$$Y_i := S_i | E_{i-1}, \dots, E_0 \quad \text{and} \quad Z_i := S_i | E_i, \dots, E_0 \quad (i \geq 1).$$

To prove the lemma it is then certainly enough to show that

$$F_{Y_i}(x) \geq F_{S_i}(x) \tag{44}$$

for all $x \in \mathbb{R}$ and $\forall i$, simply by setting $x = b_i$.

To give a proof of (44) by induction, we note first that $F_{Y_0}(x) \geq F_{S_0}(x)$ is trivial. Next, we assume that (44) holds for i and show that it holds also for $i+1$. To this end, we note first that Bayes' rule [46] implies that

$$F_{Z_i}(x) = \left\{ \begin{array}{ll} \frac{F_{Y_i}(x)}{F_{Y_i}(b_i)} & \text{if } x \leq b_i \\ 1 & \text{otherwise} \end{array} \right\} \geq F_{Y_i}(x). \tag{45}$$

The key to the proof, however, is to note is that $Y_{i+1} = Z_i + X_{i+1}$ where X_{i+1} is independent of S_j , and hence of E_j for $j \leq i$. In short, X_{i+1} is independent of Z_i . This allows to write

$$\begin{aligned}
F_{Y_{i+1}}(x) &= \mathbf{P}(Z_i + X_{i+1} < x) \\
&= \int_{-\infty}^{\infty} \int_{-\infty}^{x-x_{i+1}} f_{Z_i}(z_i) f_{X_{i+1}}(x_{i+1}) \, dz_i \, dx_{i+1} \\
&= \int_{-\infty}^{\infty} F_{Z_i}(x - x_{i+1}) f_{X_{i+1}}(x_{i+1}) \, dx_{i+1} \\
&\geq \int_{-\infty}^{\infty} F_{Y_i}(x - x_{i+1}) f_{X_{i+1}}(x_{i+1}) \, dx_{i+1} && \text{from (45)} \\
&\geq \int_{-\infty}^{\infty} F_{S_i}(x - x_{i+1}) f_{X_{i+1}}(x_{i+1}) \, dx_{i+1} && \text{from (44)} \\
&= \mathbf{P}(S_i + X_{i+1} < x) \\
&= F_{S_{i+1}}(x)
\end{aligned} \tag{46}$$

This proves the claim by induction. ◇

References

- [1] W. Leland, M. Taqqu, W. Willinger, and D. Wilson, “On the self-similar nature of Ethernet traffic (extended version),” *IEEE/ACM Trans. Networking*, pp. 1–15, 1994.
- [2] A. Erramilli, O. Narayan, and W. Willinger, “Experimental queueing analysis with long-range dependent traffic,” *IEEE/ACM Transactions on Networking*, vol. 4, pp. 209–223, April 1996.
- [3] F. Brichet, J. Roberts, A. Simonian, and D. Veitch, “Heavy traffic analysis of a fluid queue fed by a superposition of ON/OFF sources,” *COST*, vol. 242, 1994.
- [4] N. Likhanov, B. Tsybakov, and N. Georganas, “Analysis of an ATM buffer with self-similar input traffic,” *Proc. IEEE, Info com '95 (Boston 1995)*, pp. 985–992, 1995.
- [5] M. Taqqu and J. Levy, *Using renewal processes to generate LRD and high variability*. In: Progress in probability and statistics, E. Eberlein and M. Taqqu eds., vol. 11. Birkhaeuser, Boston, 1986. pp 73–89.
- [6] J. Choe and N. Shroff, “Supremum distribution of gaussian processes and queueing analysis including long-range dependence and self-similarity,” *Stochastic Models* submitted, 1997.
- [7] I. Norros, “On the use of fractional Brownian motion in the theory of connectionless networks,” *COST*, vol. 242, 1994.
- [8] N. Duffield and N. O’Connell, “Large deviations and overflow probabilities for the general single-server queue, with applications,” *Math. Proc. Cambr. Phil. Soc.*, vol. 118, pp. 363–374, 1995.
- [9] I. Norros, “Four approaches to the fractional Brownian storage,” *Fractals in Engineering*, pp. 154–169, 1997.

- [10] G. Gripenberg and I. Norros, "On the prediction of fractional Brownian motion," *J. Applied Probability*, vol. 33, pp. 400–410, 1996.
- [11] M. Crovella and A. Bestavros, "Self-similarity in World Wide Web traffic. Evidence and possible causes," in *Proceedings of SIGMETRICS '96*, May 1996.
- [12] W. Willinger, M. Taqqu, R. Sherman, and D. Wilson, "Self-similarity through high-variability: Statistical analysis of ethernet LAN traffic at the source level," *IEEE/ACM Trans. Networking (Extended Version)*, vol. 5, pp. 71–86, Feb. 1997.
- [13] N. Duffield, "Economies of scale for long-range dependent traffic in short buffers," *Telecommunication Systems*, to appear, 1998.
- [14] B. K. Ryu and A. Elwalid, "The Importance of Long-range Dependence of VBR Video Traffic in ATM Traffic Engineering: Myths and Realities," *Proc. ACM SIGCOMM Conf.*, vol. 26, no. 4, pp. 3–14, 1996.
- [15] D. P. Heyman and T. V. Lakshman, "What are the implications of long-range dependence for VBR-video traffic engineering?," *IEEE/ACM Transactions on Networking*, vol. 4, pp. 301–317, June 1996.
- [16] A. L. Neidhardt and J. L. Wang, "The concept of relevant time scales and its application to queuing analysis of self-similar traffic," in *Proc. SIGMETRICS '98/PERFORMANCE '98*, pp. 222–232, 1998.
- [17] M. Grossglauser and J.-C. Bolot, "On the relevance of long-range dependence in network traffic," *Computer-Communication-Review*, vol. 26, pp. 15–24, October 1996.
- [18] P. Flandrin, "Wavelet analysis and synthesis of fractional Brownian motion," *IEEE Trans. Inform. Theory*, vol. 38, pp. 910–916, Mar. 1992.
- [19] L. Kaplan and C.-C. Kuo, "Fractal estimation from noisy data via discrete fractional Gaussian noise (DFGN) and the Haar basis," *IEEE Trans. Signal Proc.*, vol. 41, pp. 3554–3562, Dec. 1993.
- [20] G. W. Wornell, "A Karhunen-Loève like expansion for $1/f$ processes via wavelets," *IEEE Trans. Inform. Theory*, vol. 36, pp. 859–861, Mar. 1990.
- [21] A. Tewfik and M. Kim, "Correlation structure of the discrete wavelet coefficients of fraction Brownian motion," *IEEE Trans. Inform. Theory*, vol. 38, pp. 904–909, Mar. 1992.
- [22] S. Ma and C. Ji, "Modeling video traffic in the wavelet domain," in *Proc. of 17th Annual IEEE Conf. on Comp. Comm., INFOCOM*, pp. 201–208, Mar. 1998.
- [23] L. Kaplan and C.-C. Kuo, "Extending self-similarity for fractional Brownian motion," *IEEE Trans. Signal Proc.*, vol. 42, pp. 3526–3530, Dec. 1994.
- [24] J. Roberts, U. Mocci, and J. V. (eds.), "Broadband network teletraffic," in *Lecture Notes in Computer Science, No 1155*, Springer, 1996.
- [25] S. Bates and S. McLaughlin, "The estimation of stable distribution parameters from teletraffic data," *preprint*, 1998.
- [26] R. Riedi, M. S. Crouse, V. Ribeiro, and R. G. Baraniuk, "A multifractal wavelet model with application to TCP network traffic," *IEEE Trans. Info. Theory, Special issue on multiscale statistical signal analysis and its applications*, vol. 45, pp. 992–1018, April 1999.

- [27] J. L. Véhel and R. Riedi, “Fractional Brownian motion and data traffic modeling: The other end of the spectrum,” *Fractals in Engineering*, pp. 185–202, Springer 1997.
- [28] A. C. Gilbert, W. Willinger, and A. Feldmann, “Scaling analysis of random cascades, with applications to network traffic,” *IEEE Trans. on Info. Theory, Special issue on multiscale statistical signal analysis and its applications*, April 1999.
- [29] R. H. Riedi and W. Willinger, “Toward an improved understanding of network traffic dynamics,” in *Self-similar Network Traffic and Performance Evaluation*, Wiley, June 1999.
- [30] V. Ribeiro, R. Riedi, M. S. Crouse, and R. G. Baraniuk, “Simulation of non-gaussian long-range-dependent traffic using wavelets,” *Proc. SigMetrics*, pp. 1–12, May 1999.
- [31] D. Cox, “Long-range dependence: A review,” *Statistics: An Appraisal*, pp. 55–74, 1984.
- [32] M. Taqqu, V. Teverovsky, and W. Willinger, “Estimators for long-range dependence: An empirical study,” *Fractals.*, vol. 3, pp. 785–798, 1995.
- [33] P. Abry, P. Gonçalves, and P. Flandrin, “Wavelets, spectrum analysis and $1/f$ processes,” *preprint*, 1996.
- [34] P. Abry, P. Flandrin, M. Taqqu, and D. Veitch, “Wavelets for the analysis, estimation and synthesis of scaling data,” in *Self-similar Network Traffic and Performance Evaluation*, Wiley, June 1999.
- [35] C. S. Burrus, R. A. Gopinath, and H. Guo, *Introduction to Wavelets and Wavelet Transforms: A Primer*. Prentice Hall, 1998.
- [36] I. Daubechies, *Ten Lectures on Wavelets*. New York: SIAM, 1992.
- [37] R. H. Riedi, M. S. Crouse, V. Ribeiro, and R. G. Baraniuk, “A multifractal wavelet model with application to network traffic,” *IEEE Trans. Info. Theory, (Special issue on multiscale statistical signal analysis and its applications)*, vol. 45, pp. 992–1018, April 1999. Available at www.dsp.rice.edu.
- [38] K. E. Timmerman and R. D. Nowak, “Multiscale Bayesian estimation of Poisson intensities,” in *Proc. 31st Asilomar Conf.*, (Pacific Grove, CA), Nov. 1997.
- [39] N. Johnson, S. Kotz, and N. Balakrishnan, *Continuous Univariate Distributions*, vol. 1-2. New York: John Wiley & Sons, 1994.
- [40] V. Paxson and S. Floyd, “Wide-area traffic: The failure of Poisson modeling,” *IEEE/ACM Transactions on Networking*, vol. 3, pp. 226–244, 1995.
- [41] M. Paulekar and A. M. Makowski, “Tail probabilities for a multiplexer with self-similar traffic,” *Proc. IEEE INFOCOM*, pp. 1452–1459, 1996.
- [42] B. V. Rao, K. R. Krishnan, and D. P. Heyman, “Performance of Finite-Buffer Queues under Traffic with Long-Range Dependence,” *Proc. IEEE GLOBECOM*, vol. 1, pp. 607–611, November 1996.
- [43] D. V. Lindley, “The theory of queues with a single server,” *Proceedings of the Cambridge Philosophical Society*, vol. 48, pp. 277–289, 1952.
- [44] D.-Y. Fan, “The distribution of the product of independent beta variables,” *Commun. Statist.-Theory Meth.*, vol. 20, no. 12, pp. 4043–4052, 1991.

- [45] R. H. Riedi, "Multifractal processes," *Technical Report, ECE Dept. Rice Univ., TR 99-06*, submitted for publication 1999.
- [46] A. Papoulis, *Probability, Random Variables, and Stochastic Processes*. New York: McGraw-Hill, 1991.

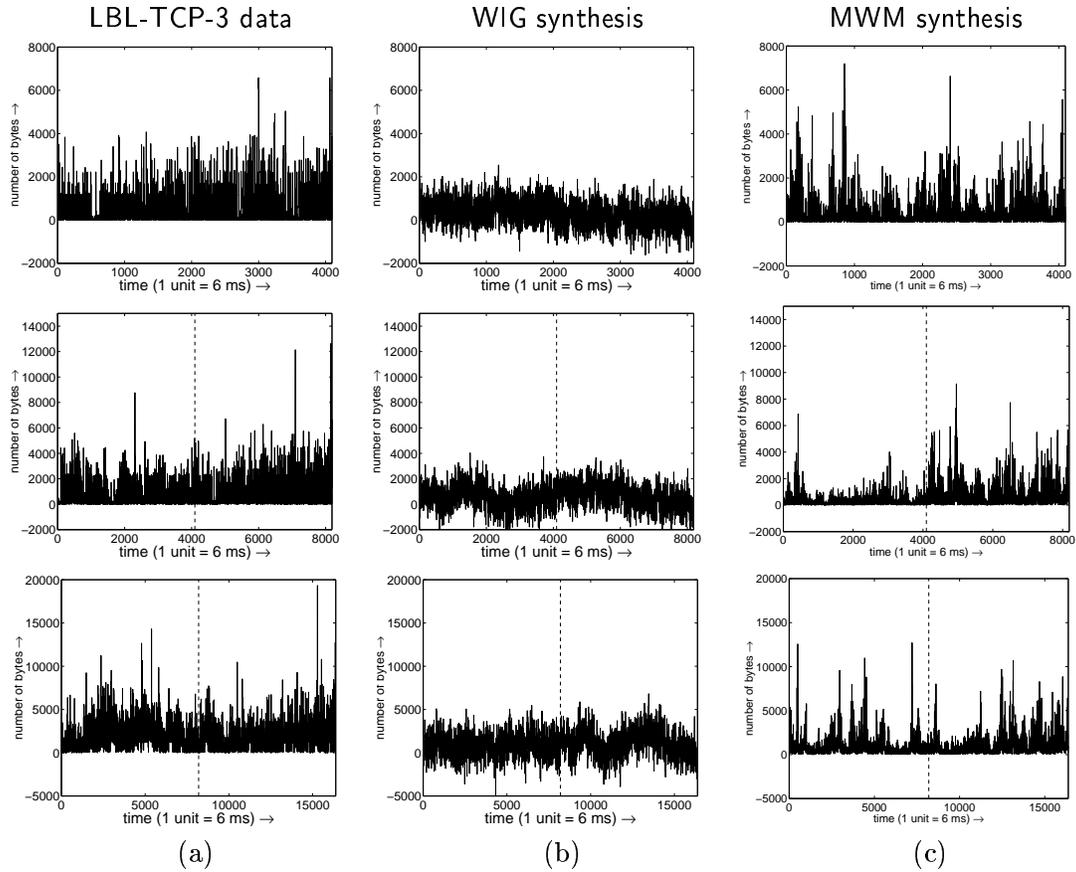


Figure 1: Bytes-per-time arrival process at different aggregation levels for (a) wide-area TCP traffic at the Lawrence Berkeley Laboratory (trace LBL-TCP-3) [40], (b) one realization of the state-of-the-art wavelet-domain independent Gaussian (WIG) model [22], and (c) one realization of the multifractal wavelet model (MWM) synthesis. The top, middle and bottom plots correspond to bytes arriving in intervals of 6ms, 12ms and 24ms respectively. The top and middle plots correspond to the second half of the middle and bottom plots, respectively, as indicated by the vertical dotted lines. The MWM traces closely resemble the real data, while the WIG traces (with their large number of negative values) do not.

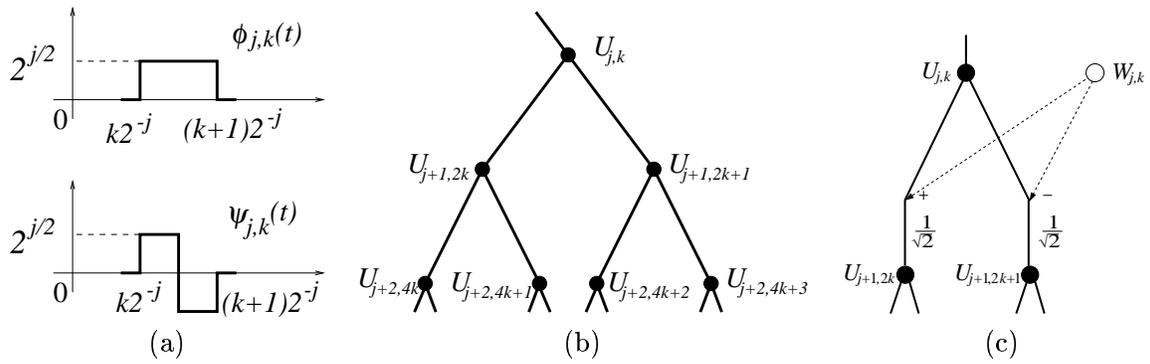


Figure 2: (a) The Haar scaling and wavelet functions $\phi_{j,k}(t)$ and $\psi_{j,k}(t)$. (b) Binary tree of scaling coefficients or local mean values of the signal. Nodes at each horizontal level in the tree are coarse representation of the signal with lower levels corresponding to finer resolutions of the signal. (c) Recursive scheme for calculating the Haar scaling coefficients $U_{j+1,2k}$ and $U_{j+1,2k+1}$ at scale $j + 1$ as sums and differences (normalized by $1/\sqrt{2}$) of the scaling and wavelet coefficients $U_{j,k}$ and $W_{j,k}$ at scale j . For the WIG model, the $W_{j,k}$'s are mutually independent and identically distributed within scale according to $W_{j,k} \sim N(0, \sigma_j^2)$.

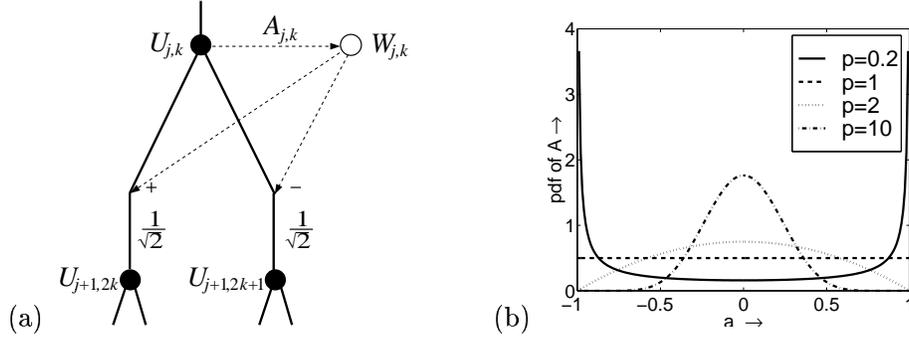


Figure 3: (a) Multifractal wavelet model (MWM) construction: At scale j , generate the multiplier $A_{j,k} \sim \beta_{-1,1}(p_j, p_j)$, and then form the wavelet coefficient as the product $W_{j,k} = A_{j,k}U_{j,k}$. At scale $j + 1$, form the scaling coefficients in the same manner as the WIG model in Figure 2(c). (b) Probability density function of a $\beta_{-1,1}(p, p)$ random variable A . For $p = 0.2$, $\beta_{-1,1}(p, p)$ resembles a binomial distribution, and for $p = 1$ it has a uniform density. For $p > 1$ the density is close to a truncated Gaussian density with increasing resemblance as p increases.

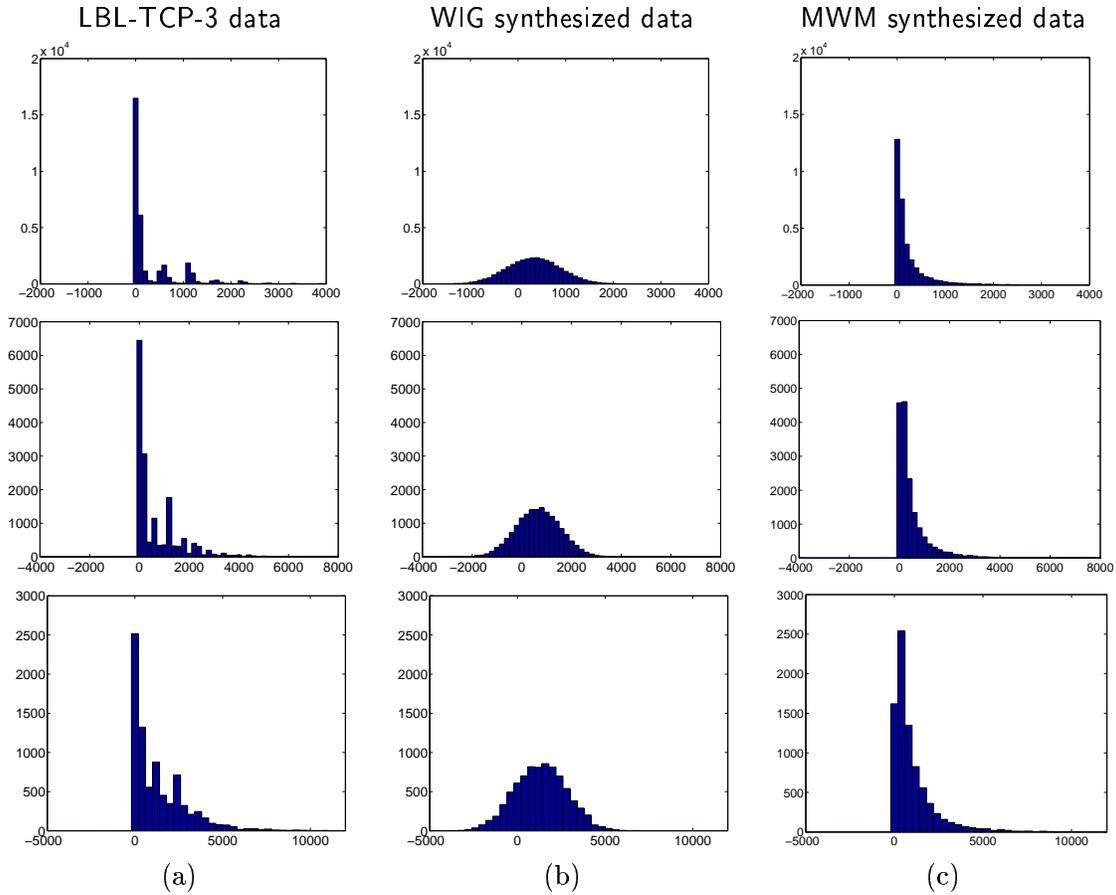


Figure 4: Histograms of the bytes-per-times process at different aggregation levels for (a) wide-area TCP traffic at the Lawrence Berkeley Laboratory (trace LBL-TCP-3) [40], (b) one realization of the WIG model, and (c) one realization of the MWM synthesis. The top, middle and bottom plots correspond to bytes arriving in intervals of 6ms, 12ms and 24ms respectively. Note the large probability mass over negative values for the WIG model.

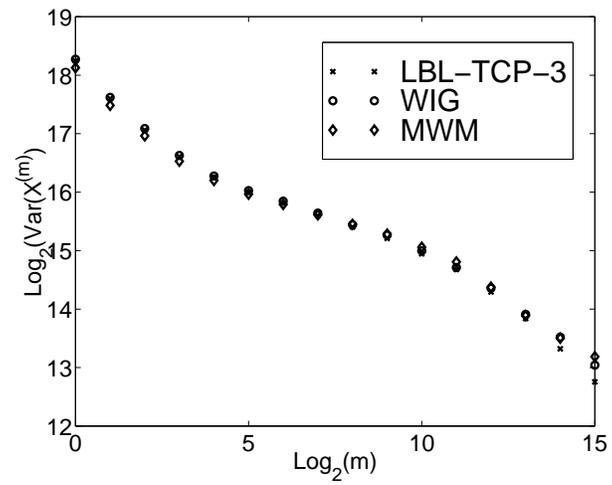


Figure 5: Variance-time plot of a real data trace (LBL-TCP-3) and one realization of each, the WIG and the MWM synthesis. Both, the MWM and WIG model, capture the correlation structure of the real data.

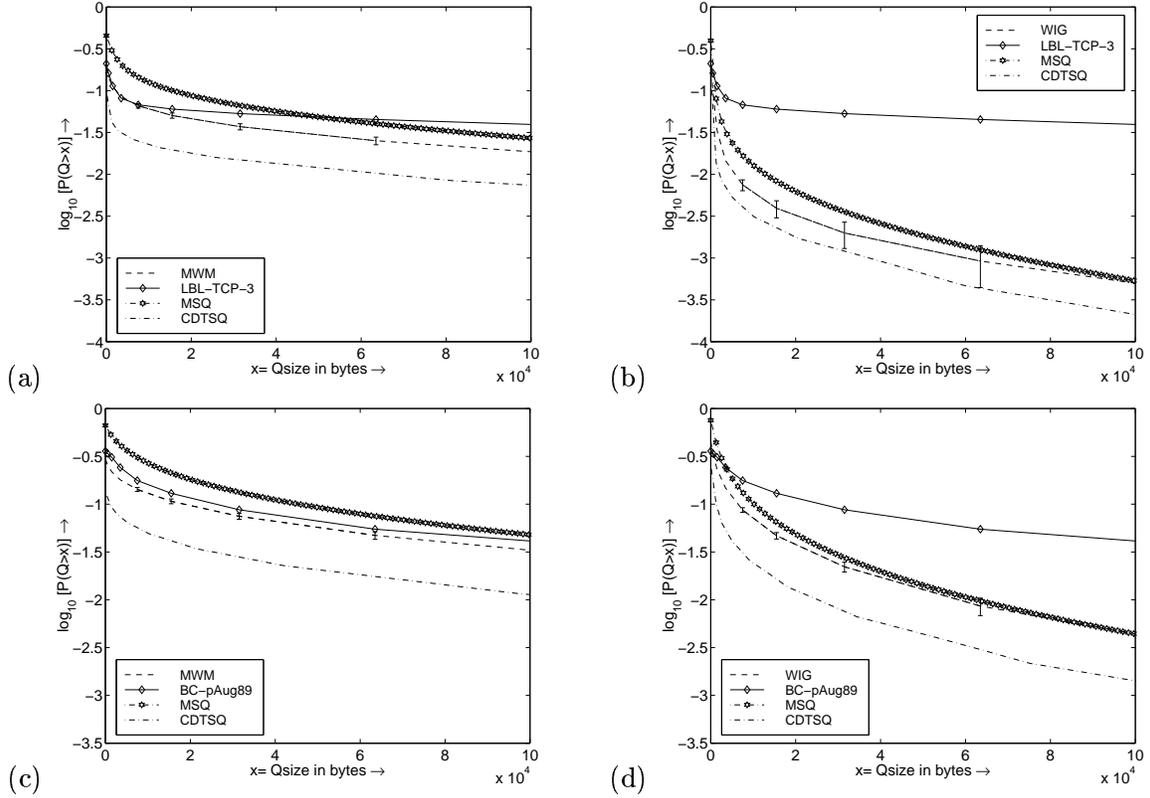


Figure 6: Comparison of the queuing performance of real data traces with those of synthetic WIG and MWM traces. In (a), we observe that the MWM synthesis matches the queuing behavior of the LBL-TCP-3 data closely, while in (b) the WIG synthesis does not. In (c) and (d), we observe a similar behavior with the BC-pAug89 data. We also observe that the MSQ is a close approximation to the empirical queuing behavior for both synthetic traffic loads, the WIG and MWM and that it is closer than the CDTSQ.

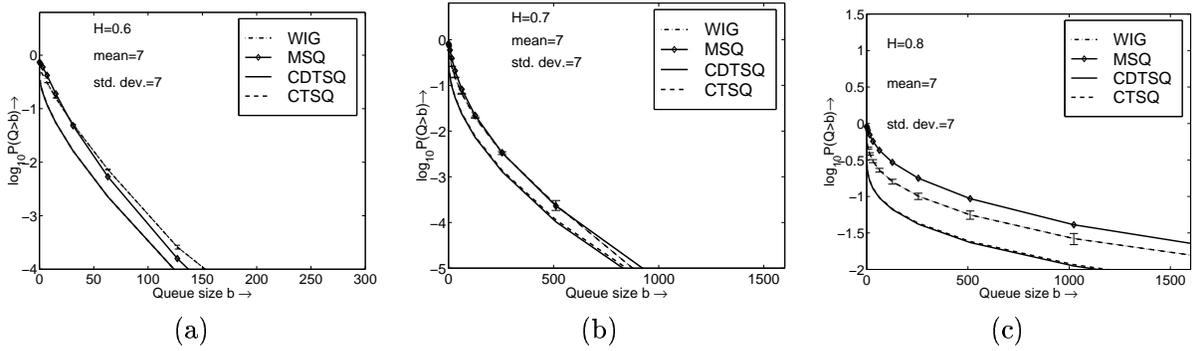


Figure 7: Justification for limiting the queuing analysis to dyadic time scales through comparison of the critical time scale queue (CTSQ) and the critical dyadic time scale queue (CDTSQ) for the WIG. Experiments used synthetic WIG traces corresponding to an fGn correlation structure for different values of Hurst parameter H . In (a) $H = 0.6$, in (b) $H = 0.7$ and in (c) $H = 0.8$. In all cases, the mean, standard deviation and link capacity were 7, 7 and 10 units respectively. Observe that in all cases the CDTSQ and CTSQ are almost identical.

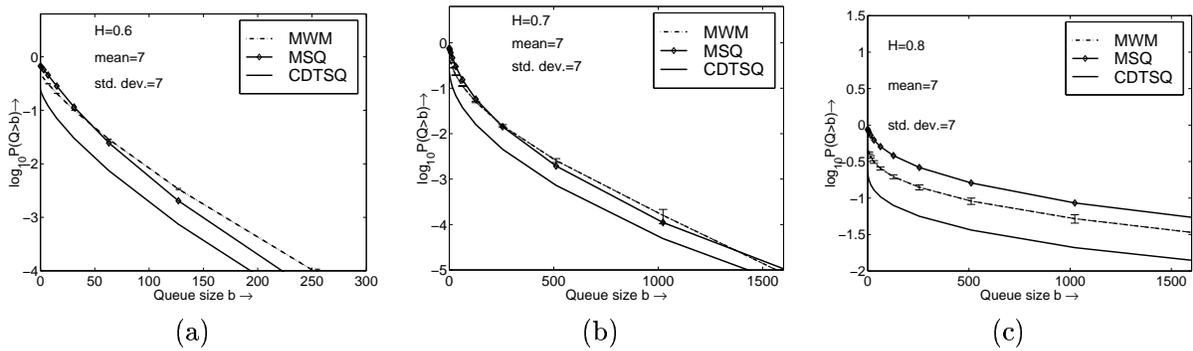


Figure 8: Validation of the theoretical formula (28) for the tail queue probability of the MWM and WIG models. These experiments used synthetic MWM traces corresponding to an fGn correlation structure for different values of Hurst parameter H . In (a) $H = 0.6$, in (b) $H = 0.7$ and in (c) $H = 0.8$. In all cases, mean, standard deviation and link capacity were 7, 7 and 10 units respectively. Observe that in all cases the MSQ formula gives a good approximation to the empirical queuing behavior.

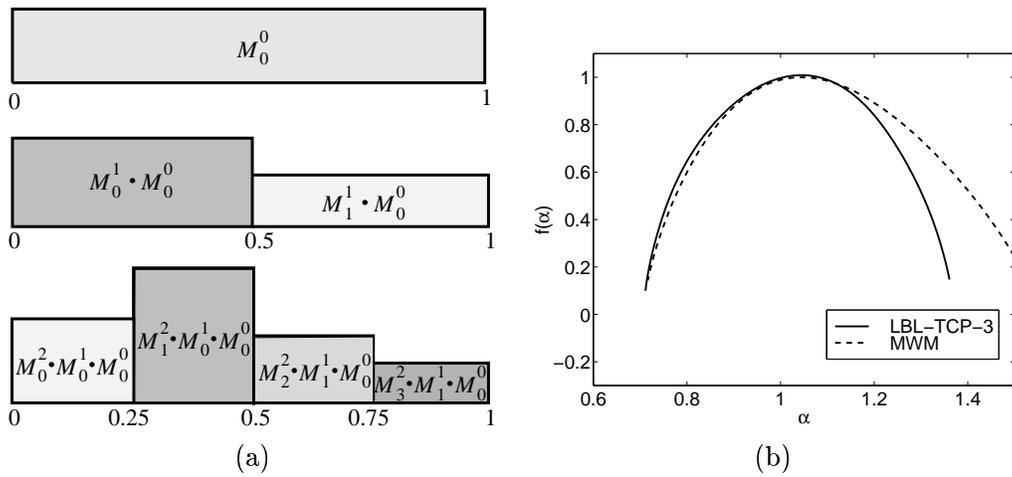


Figure 9: (a): The MWM translates immediately into a multiplicative cascade in the time domain (cf. (34)). (b) Multifractal spectra of the LBL-TCP-3 data and one realization of the MWM synthesis. The MWM spectrum matches that of the real data closely except for large values of α or small values of the signal.