# Analyzing the relevant time scales in a network of queues

A. Nogueira[*], R. Valadas[†]
Institute of Telecommunications - University of Aveiro
Campus Universitário Santiago
3810-193 Aveiro, Portugal

## ABSTRACT

Network traffic processes can exhibit properties of self-similarity and long-range dependence, i.e., correlations over a wide range of time scales. However, as already shown by several authors for the case of a single queue, the second-order behavior at time scales beyond the so-called correlation horizon or critical time scale does not significantly affect network performance. In this work, we extend previous studies to the case of a network with two queuing stages, using discrete event simulation. Results show that the second stage provokes a decrease in the correlation horizon, meaning that the range of time scales that need to be considered for accurate network performance evaluation is lower than predicted by a single stage model. We also used simulation to evaluate the single queue model. In this case, the estimated correlation horizon values are compared with those predicted by a formula derived by Grossglauser and Bolot, which presumes the approximation of the input data by a traffic model that enables to control the autocorrelation function independently of first-order statistics. Results indicate that although the correlation horizon increases linearly with the buffer size in both methods, the simulation ones predict a lower increase rate.

Keywords: self-similarity, long-range dependence, wavelets, time-scales, correlation horizon, critical time scale, packet loss ratio.

## 1. INTRODUCTION

There is wide experimental evidence that network traffic processes exhibit properties of self-similarity and long-range dependence (LRD)[2-4]. For the case of a single queue, it was shown that second-order behavior at the time scales beyond the correlation horizon (CH) or critical time scale (CTS) does not significantly affect network performance[1,5].

CH and CTS are two different terms representing essentially the same concept: the time scale or lag that separates relevant and irrelevant correlation with respect to the performance measure of interest, which was assumed to be the packet loss ratio (PLR). The CH depends on the correlation structure of the input traffic and on the system under study. In real networks, the buffers have finite length. A finite length buffer "forgets" about the past as soon as it is either empty or full. Thus, while correlation on all time scales has an impact on the performance of an infinite queue, only the correlation up to the CH has an effect in a finite buffer queue. It was shown in Ref. 5 that the CH is finite, has a small value for small buffer and is a non-decreasing function of the buffer size.

It is known that processes with the same correlation structure can generate vastly different queuing behavior. In fact, performance evaluation depends not only on the time scales relevant to the system under study and on the correlation structure of the source traffic but also on other characteristics like the marginal distribution of the arrival rate process[1,6]. So, it is important to capture and study the relative influence of each one of these factors in the overall performance. However, our goal in this paper is only to study the influence of the correlation structure of the source traffic and the characteristics of the queuing system in the overall queuing performance, leaving other parameters such as the marginal distribution of the input traffic fixed.

In order to calculate the amount of correlation that needs to be taken into account for performance evaluation and to analyze its behavior along the various stages of a queuing network, we consider a realistic scenario consisting of buffers (Figure 1) with adjustable parameters (queue length, link rate) and an input traffic stream exhibiting the LRD property. We perform several tasks: (i) analysis of the LRD characteristics of the input traffic stream at each node, using the autocorrelation

[*] e-mail: nogueira@av.it.pt; phone: +351 234377900; fax: +351 234377901

[†] e-mail: rv@det.ua.pt; phone: +351 234377900; fax: +351 234377901

function and the Wavelet estimator[7]; (ii) estimation of the CH in each network node, via simulation; (iii) comparison of the CH values corresponding to different stages; (iv) for the single buffer case, comparison of the CH values obtained using our simulation approach and the method proposed in Ref. 1.
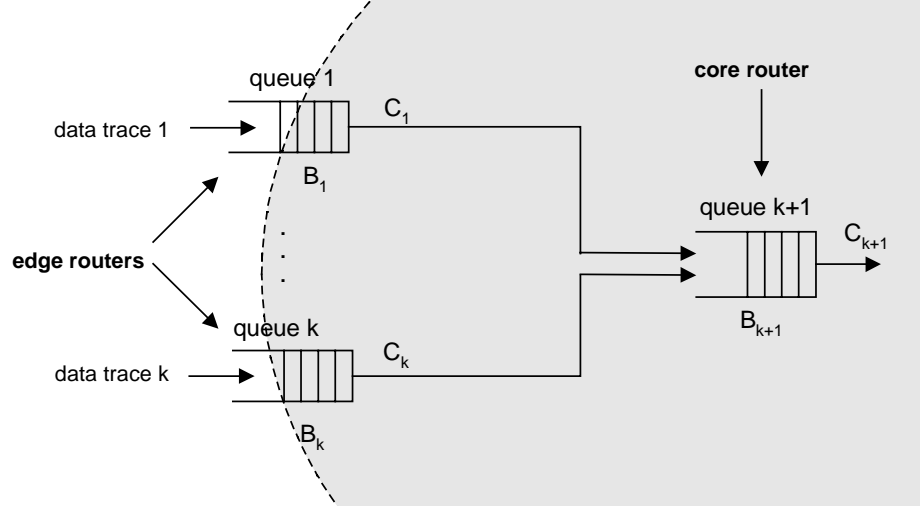


Figure 1: Network of finite queues used in the study of the correlation horizon

The network that constitutes the basis for our simulation studies tries to replicate a realistic scenario consisting of a network of switching elements, e.g., routers (Figure 1). The buffers of stage 1 belong to edge routers, located at the entry points of the network. In our case, for the sake of simplicity, we restrict our scenario to two input queues only ( $k = 2$ ). The buffer of stage 2 belongs to a core router, located at the core of the network and is typically fed by the superposition of several traffic streams, which are output flows of other queues.

The rest of the paper is structured as follows: in section 2, we present an outline of the methodology used in this study; section 3 analyzes the characteristics of the input trace, including the study of its LRD behavior using the wavelet estimator tool; in section 4, the main results are presented and discussed; finally, in section 5, we draw the main conclusions.

## 2. OVERVIEW OF THE METHODOLOGY

Our study of the correlation horizon is based on discrete event simulation: using the August Bellcore data trace[2], which will be designated in the rest of this paper by pAug, as the input traffic stream 1, we control the "extent" of its correlation using a shuffling technique. Data trace 2 in Figure 1 represents the traffic flow of another set of services and is simply a shifted version of the original pAug data trace. We will denote this trace by pAugShift.

The external shuffling procedure, described in Ref. 8, eliminates correlation in the input process beyond a certain lag. In this procedure, a time series representing a realization of a process is divided up into blocks and the blocks are shuffled. However, the structure of the time series inside a block remains unchanged. In this way, external shuffling removes correlation from the time series beyond a lag equal to the length of a block.

The main parameters of the queuing systems of both stages are selected such that the PLR values have the same order of magnitude. For stage 1, we chose an output link capacity $C_1$ for both buffers, according to the average arrival rate $\lambda_i$ of the input data traces. For the stage 2 buffer, the average arrival rate is $\lambda_o < \lambda_i$. In order to have the same utilization ratio in both stages, we selected a capacity $C_3$ slightly lower than $2C_1$. In our simulations, we considered $C_1 = 1.4$ Mbit/s and $C_3 = 2.1$ Mbit/s.

Applying the pAug data trace to queue n° 1, we calculate its CH value by determining the PLR values for varying block lengths (defined according to the shuffling method) and normalized buffer sizes, $B / C$. The normalized buffer size represents the time it takes to empty a buffer of length $B$ packets at an output rate of $C$ packets/s.

In a second step, the output flow of queue nº 1 is applied to queue nº 2 together with any other data traces coming from other queues, which are represented in our simulations by the pAugShift trace. One again, the CH value is calculated through the determination of the PLR curve for varying block lengths and normalized buffer sizes. The CH values calculated in both stages of the network are then compared, taking into account the parameters of the buffer system.

In order to evaluate the accuracy of the estimated CH values in the first stage queues, we approximate the input data trace by the traffic model proposed in Ref. 1 which enables us to examine the impact on the performance measure of interest (PLR, in this case) of considering only the relevant time scales of the traffic process. The parameters of this N state traffic model are estimated from the empirical trace using a fitting procedure that will be described at the end of this section. The traffic model in each state is a process with constant rate and the duration of each state is represented by a truncated Pareto distribution, with cumulative distribution function given by:

$$F_T(t) = \begin{cases} \left(\dfrac{t+\theta}{\theta}\right)^{-\alpha} & \text{if } t < T_C \\ 0 & \text{otherwise} \end{cases}$$

(1)

where $1 < \alpha < 2$ represents the exponential decay. The parameter $T_C$ is called the cutoff lag and, since the duration of a state cannot exceed $T_C$ and since the rates in consecutive intervals are independent, there is no correlation in the rate process beyond lag $T_C$. Therefore, this parameter represents the correlation horizon concept.

This traffic model enables to control individually the correlation structure of the data trace, keeping other characteristics of the process (such as the marginal distribution) unchanged. The fitting algorithm used to estimate its parameters from the empirical trace is represented schematically in Figure 10, at the end of this paper, and is based on the fitting approach presented in Ref. 1. The vectors $\Pi$ and $\Lambda$, representing the marginal rate distribution and the rate matrix, respectively, are estimated from the histogram of the Number of Arrivals (NAF) function. In the construction of this histogram, we assumed 50 bins in all our experiments, as done in Ref. 1. In order to completely define the truncated Pareto distribution of equation (1), one must estimate the parameters $T_C$, $\theta$ and $\alpha$. In order to estimate $\theta$, we first calculate the average number of consecutive samples in the trace that fall within the same histogram bin, $E[T_n]$. Then, $\theta$ is calculated assuming that the mean interval duration, which is given by

$$E[T_n] = \frac{\theta}{\alpha - 1}\left[1 - \left(\frac{T_c}{\theta} + 1\right)^{1-\alpha}\right]$$

(2)

matches the empirical mean. We calculate the value of $\theta$ for each bin using a minimization procedure.

The calculation and analysis of the cutoff lag or correlation horizon, $T_C$, which represents the main focus of this study, is performed in two different ways. First, using the results presented in Ref. 1 where an explicit expression is given

$$T_{CH} = \frac{B}{2\sqrt{2}\sigma_T\sigma_\lambda erf^{-1}(p)}(\mu + \beta\sigma_T)$$

(3)

which includes first and second order statistics of the interarrival time and arrival rate distribution functions. Second, using our simulation approach. Finally, the exponential decay $\alpha$ is calculated from the Hurst parameter, which can be estimated using various estimation techniques, like for example the variance time-plot.

Once we have completely defined the Pareto distribution, describing the duration of each state, and estimated vectors $\Pi$ and $\Lambda$, the traffic model is completely estimated from the empirical trace and simulated data traces can be generated using discrete event simulation.

## 3. CHARACTERISTICS OF THE INPUT TRACE

Trace pAug has a total of $10^6$ values (packet arrival instants and their respective packet sizes), the average packet size is 434.29 bytes and the average interarrival time value is 0.0031 s. The Hurst parameter, estimated through the variance-time

plot, is $0.8745$ which gives and $\alpha$ value of $1.2509$. In the conversion step from the arrival instant to the number of arrivals formats, the sample interval was set to $10$ ms.

We analyze the presence of LRD behavior in the input traces using the scaling analysis method described in Ref. 7. This method resorts to the so-called Logscale Diagram which consists in the graph of $y_j$ against $j$, together with confidence intervals about the $y_j$, where $y_j$ is a function of the wavelet discrete transform coefficients at scale $j$. Traffic is said to be LRD if, within the limits of the confidence intervals, the $y_j$ fall on a straight line, in a range of scales from some initial value $j_1$ up to the largest one present in data. The results of this scaling analysis are present in Figure 2 for trace pAug.
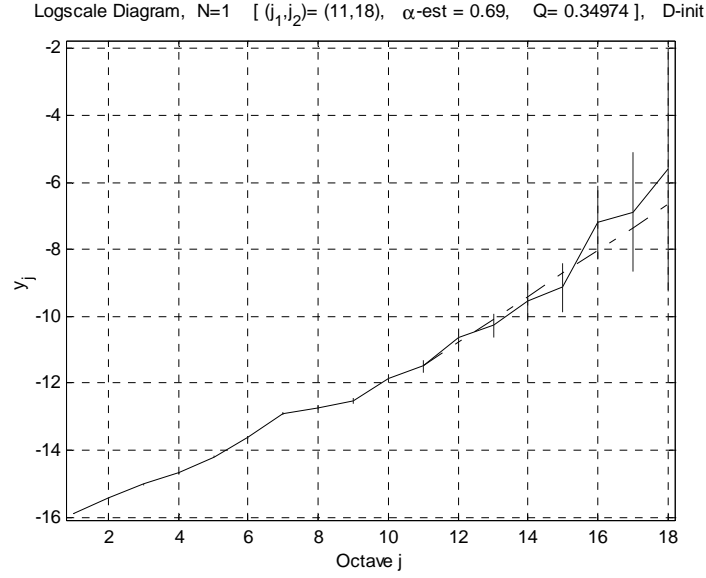


Figure 2: Scaling analysis for trace pAug.

The analysis of the autocovariance function, represented in Figure 3, lead us to suspect that the trace has a LRD behavior, due to the slow decay for large time lags. This is confirmed by the scaling analysis, since the $y_j$ values are aligned between a medium octave, 11, and octave 18, the highest octave present in data.
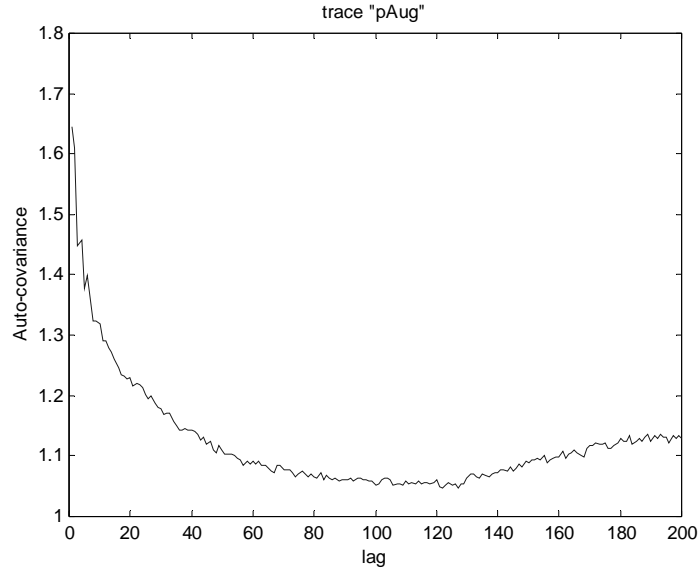


Figure 3: Autocovariance function of the interarrival times for trace pAug

## 4. RESULTS AND DISCUSSION

Applying the pAug data trace to queue nº 1, we simulated the *PLR vs block size* function, for different normalized buffer sizes. The results are plotted in Figure 4. We used normalized buffer sizes of up to 2 seconds, which are typical values for currently available switches. For example, an ATM switch with a buffer of 7000 cells will produce a maximum delay in the buffer approximately equal to 2 $s$, considering a T1 link at 1.5 Mbit/s. The selection of the output link capacity depends on the utilization ratio in each queue and must result in PLR values that span over several orders of magnitude. For the first stage, the capacity was set to 1.4 Mbit/s.
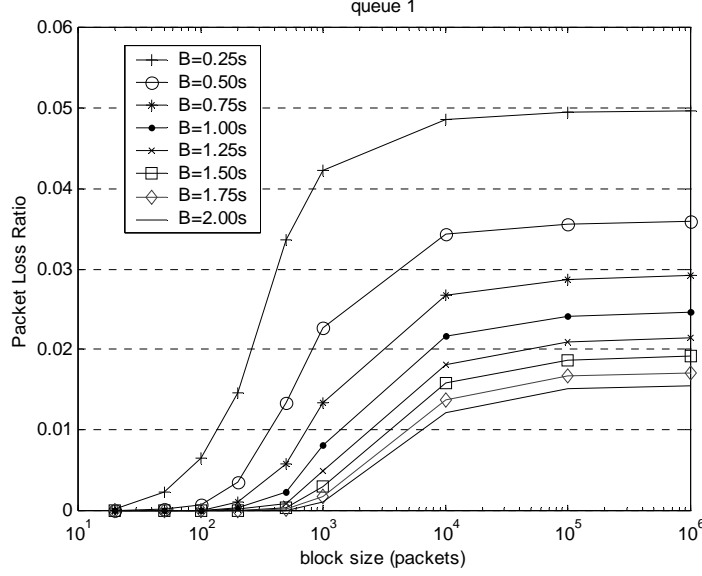


Figure 4: Packet loss ratio vs block size in the first stage, for different normalized buffer sizes, *B*.

The CH, being the time scale that separates relevant and irrelevant correlation with respect to PLR, corresponds to the block size value beyond which PLR does not change significantly. In our case, we have assumed that the CH corresponds to the block size beyond which PLR has a value less than or equal 5% of the final value (which corresponds to the unshuffled trace). This value represents the tolerance in the PLR value for the estimation of the CH.

To determine CH we interpolate linearly the *PLR vs block size* curve, and try to minimize the error by performing additional simulations around the threshold block size value (the one that corresponds to the CH). Considering a tolerance of *x*% for the PLR, the threshold block size value is given by

$$bs = \frac{(1-x)PLR_f(BS_2 - BS_1) - PLR_1 BS_2 + PLR_2 BS_1}{PLR_2 - PLR_1} \tag{4}$$

where $PLR_f$ is the PLR corresponding to the unshuffled trace, $PLR_2$ and $PLR_1$ are the PLR values for the higher and lower endpoints of the segment that includes $(1-x)PLR_f$, respectively, and $BS_2$ and $BS_1$ are the block size values corresponding to the higher and lower endpoints of the segment.

The correlation horizon values calculated for the first queue are represented in the higher curve of Figure 5.
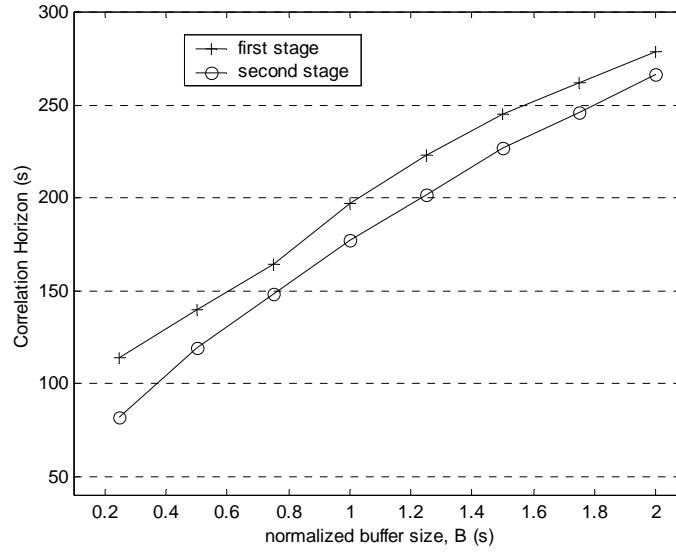
Figure 5: Comparison of the CH values calculated for both stages of the network ( $C_1 = 1.4$ Mbit/s and $C_3 = 2.1$ Mbit/s)

In a second step, the output flows of queues nº 1 and nº 2 are applied to the second stage queue. Again, we calculated the *PLR vs block size* function, for different normalized buffer sizes. The results are plotted in Figure 6, for an output link capacity of 2.1 Mbit/s. In order to evaluate the effect of the utilization ratio in the CH values, we plot in Figure 7 the *PLR vs block size* function, for different normalized buffer sizes, but for an output link capacity of 1.4 Mbit/s. In this case, the PLR values are obviously higher and they do not span over several orders of magnitude as happened in the first case. However, the differences in terms of CH values are not significant.
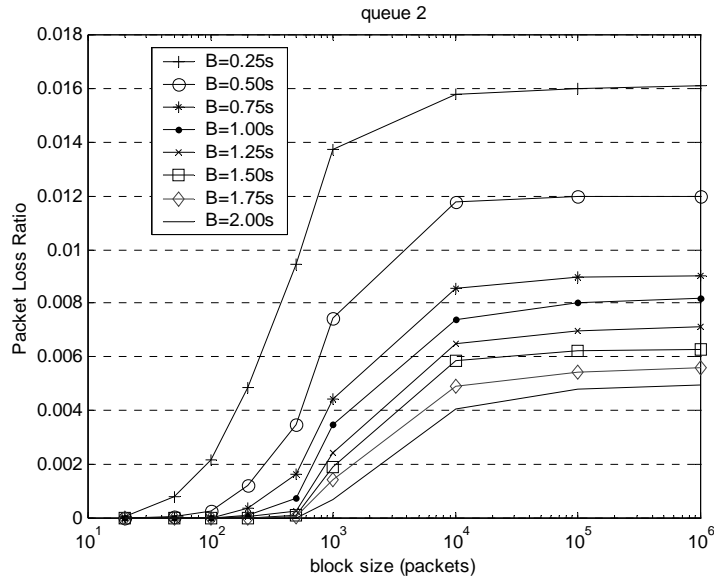


Figure 6: Packet loss ratio in the second stage as a function of the block size and for different normalized buffer sizes. The output link capacity was set to $2.1$ Mbit/s
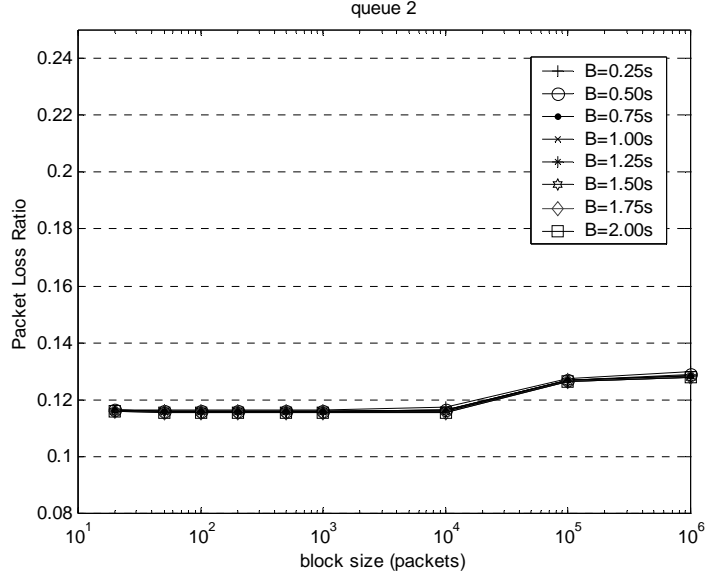
Figure 7: Packet loss ratio in the second stage as a function of the block size and for different normalized buffer sizes. The output link capacity was set to $1.4$ Mbit/s

The new CH values, corresponding to $C_3 = 2.1$ Mbit/s, are calculated using the same procedure described above and are represented in the lower curve of Figure 5. As we can see, the CH values for stage 2 are lower than the corresponding values for stage 1, which means that in the second stage of the network the relevant correlation (in terms of PLR) ends at lower lags.

From Figure 5 we also see that the CH varies almost linearly with the normalized buffer size, for both stages, as already concluded in Ref. 1. In fact, the curves obtained are not perfectly linear, but we think that this slight deviation is caused by the error introduced in our interpolation procedure described above.

The work in Ref. 1 also predicts that the slope of the *CH vs normalized buffer size* is one. This can also be seen in equation 3. However our simulation results predict a lower slope, as shown in Figure 8. Both methods lead to similar CH values only for small normalized buffer sizes, and the deviation between the two approaches increases with the normalized buffer size. For $B = 0.50$ s, for example, the difference between both values is about $48\%$. The effect, in terms of PLR, of those differences in the estimated CH values are illustrated in Figure 9.

We can conclude that expression (3) used in Ref. 1 to calculate the CH provides a pessimistic estimation of the relevant correlation. That expression was derived based on the assumption that the *buffer occupancy plus sum of the excess work in n consecutive intervals* was approximately normally distributed. Our simulation results show that the real relevant correlation that needs to be taken into account for PLR prediction is much smaller than that, making our modeling tasks of the autocorrelation function easier to perform. This in fact enlarges the set of traffic models that can be used to match long-range dependence characteristics.

## 5. CONCLUSIONS

In this paper we evaluated the amount of correlation that needs to be taken into account for performance evaluation and analyzed its behavior along both stages of a two stage queuing network. In order to do that, we (i) presented a discrete-event simulation methodology for the estimation of the correlation horizon values in each stage of the network; (ii) compared the results obtained for both stages; and (iii) compared the results obtained for the single stage case with the ones presented by Grossglauser and Bolot.

Results have shown that the CH values on stage 2 are lower than the corresponding values for stage 1, which means that in those stages the relevant correlation (in terms of packet loss ratio) ends at lower lags. For the single buffer case, Grossglauser and Bolot have concluded that when keeping the interarrival time and arrival rate distribution functions
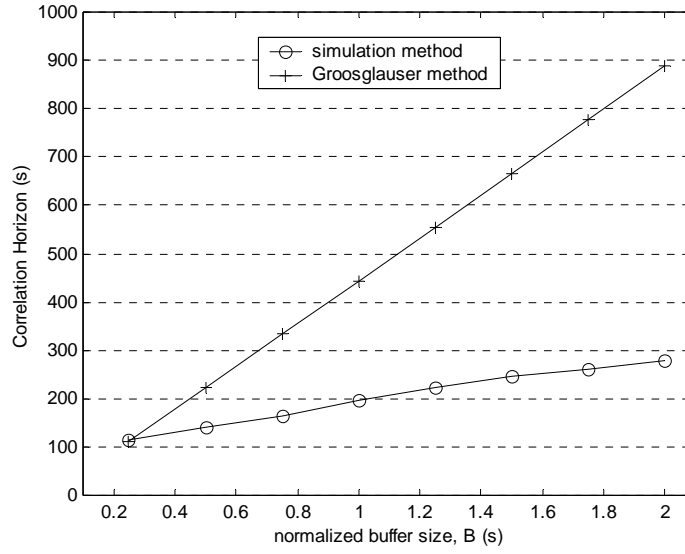
Figure 8: Comparison of the CH values, for stage 1, calculated using the simulation method and the method described in[1] (designated here by Grossglauser method)
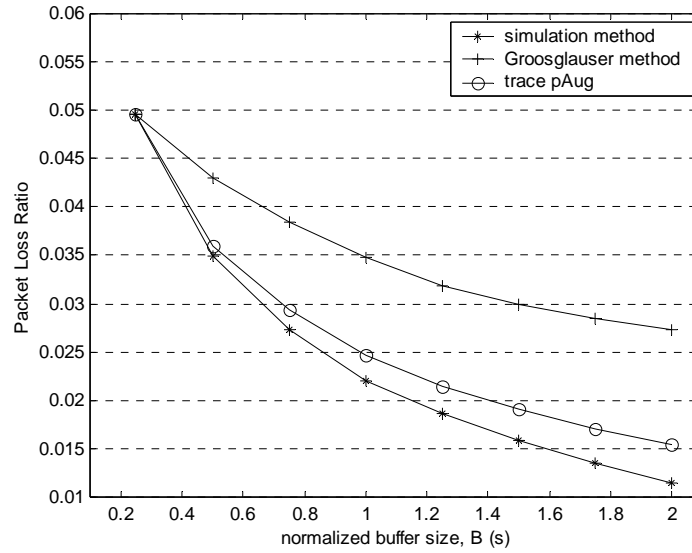


Figure 9: Comparison of the effect, in terms of PLR, of the CH values calculated using the simulation method and the method described in Ref. 1

unchanged the CH varies linearly with the normalized buffer size with a slope equal to 1. However, our simulation results predicted a lower slope. The deviation between the values estimated by both approaches increases with the normalized buffer size.

The results obtained can be applied in network dimensioning and call admission control procedures, enabling the dimensioning of each network node buffer and output link capacity in order to optimize the overall network performance.

We have shown in our study that the relevant correlation is limited to a time scale smaller than or equal to the correlation horizon. This fact enlarges the set of traffic models that can be used to match long-range dependence characteristics. The choice can be based on analytic tractability, on ease of parameter identification and estimation from empirical traces, or any other criteria. Grossglauser and Bolot use the truncated Pareto model because it is a parsimonious model, which provides a

simple way to control its correlation structure, and because it is self-similar when $T_C$ is set to infinity. However, Markovian models would be another possible choice since they can capture correlations up to a given value CH. Several studies have used Markov models to approximate traffic sources with LRD but the resulting Markov models are complex multi-state models that do not follow the principle of parsimonious modeling because every state added to such a model also adds several free parameters. Nonetheless, it is possible to reduce the impact of the growing number of parameters in multi-state Markov models modeling each time scale by a different state of the model.

## ACKOWLEDGMENTS

## REFERENCES

1.  M. Grossglauser and J. C. Bolot, "On the Relevance of Long-Range Dependence in Network Traffic", *IEEE/ACM Trans. On Networking* **7**, October 1999.
2.  W. E. Leland, M. S. Taqqu, W. Willinger and D. V. Wilson, "On the Self-Similar Nature of Ethernet Traffic", *IEEE/ACM Trans. On Networking* **2**, pp. 1-15, February 1994.
3.  V. Paxson and S. Floyd, "Wide Area Traffic: The Failure of Poisson Modelling", *IEEE/ACM Trans. On Networking* **3**, pp. 226-244, June 1995.
4.  J. Beran, R. Sherman and W. Willinger, "Long Range Dependence in Variable Bit Rate Video Traffic", *IEEE Trans. On Communications* **43**, pp. 1566-1579, February 1995.
5.  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'96*, Stanford, CA, August 1996.
6.  A. Anderson and B. Nielsen, " A Markovian Approach for Modeling Packet Traffic with Long-Range Dependence", *IEEE Journal on Selected Areas in Communications* **16**, pp. 719-732, June 1998.
7.  D. Veitch and P. Abry, "A wavelet based joint estimator for the parameters of LRD", *Special issue on Multiscale Statistical Signal Analysis and its Applications - IEEE Trans. Inform. Theory* **45**, April 1999.
8.  A. Erramilli, O. Narayan and W. Willinger, "Experimental Queueing Analysis with Long-Range Dependent Packet Traffic", *IEEE/ACM Trans. On Networking* **4**, April 1996.
9.  S. Q. Li and C. L. Hwang, "Queue Response to Input Correlation Function: Continuous Spectral Analysis", *IEEE/ACM Trans. On Networking* **1**, pp. 678-693, March 1993.
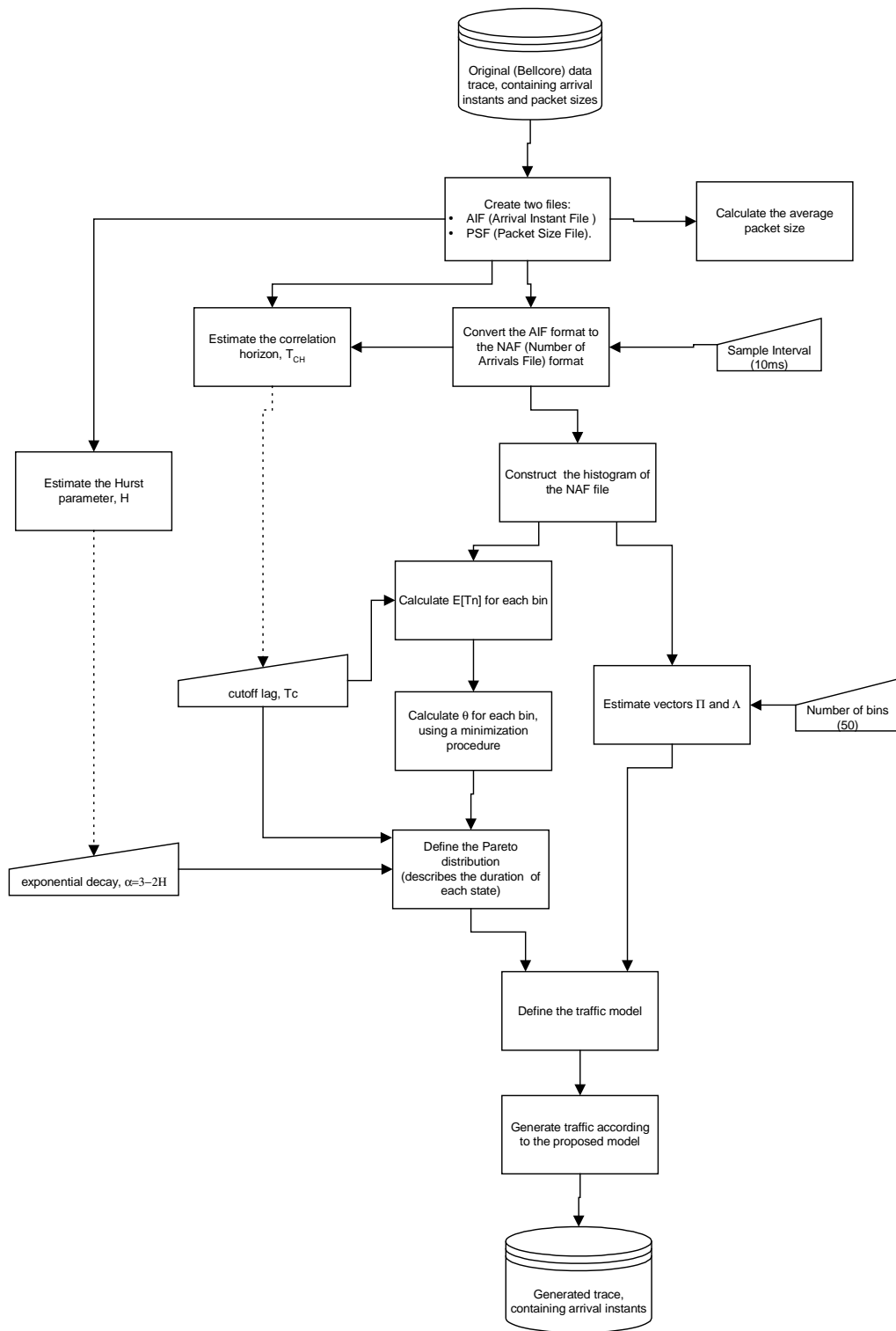
Figure 10: Fitting algorithm for the traffic model proposed to approximate data traces taking into account only their relevant correlation