# Time Series Models for Internet Data Traffic

Chun You and Kavitha Chandra
*Center for Advanced Computation and Telecommunications*
*Department of Electrical and Computer Engineering*
*University of Massachusetts Lowell*
*Lowell, MA 01854*
*E-mail: cyou@cs.uml.edu,kavitha_chandra@uml.edu*

## Abstract

*A statistical analysis of Internet traffic measurements from a campus site is carried out to examine the influence of the constituent protocols and applications on the characteristics of the aggregate stream and on packet loss statistics. While TCP remains the dominant traffic protocol through all hours of the day, a mixture of both well-known (http, ftp, nntp and smtp) and less known applications contribute significant portions to the TCP traffic mix. Statistical tests show that the aggregate TCP packet arrival process exhibits both non-stationary and nonlinear features. By filtering a subset of the applications found to exhibit non-stationary features from the aggregate process, a stationary traffic stream is derived. This filtered traffic process is modeled using nonlinear threshold autoregressive processes. The traffic model is shown to provide good agreement with the measurement trace in the packet loss statistics. The proposed parametric model allows the design of traffic shapers and provides a simple and accurate approach for simulating Internet data traffic patterns.*

## 1. Introduction

Accurate modeling of the offered traffic load is the first step towards the goal of optimizing resource allocation such that provision of services complies with the quality of service (QoS) constraints while maintaining maximum network utilization. Traffic modeling is also necessary for traffic forecasting and engineering future network capacity. Traffic measurement studies have amply demonstrated the complexity inherent in data traffic patterns [1,2,3]. Data packet arrivals are found to be correlated over both short and long-time scales. These features generally result from the arrival of bursts of packets of comparable size, often leading to high instantaneous arrival rates. Traffic modeling studies have attempted to capture the long range dependence features through self-similar fractal processes. However these models require that the underlying traffic process exhibit

properties of time-stationarity and linearity in its arrival rates. The stationarity assumption may not be easily made if one considers that both the number of sources generating traffic and the application mix change significantly with time. In addition, the statistical analysis is typically carried out on aggregated traffic measurements in which different applications and protocols are merged into a common traffic stream for analysis. To what degree this level of aggregation impacts the statistical features of traffic has not been investigated, to the best of our knowledge. It is shown in this paper that applications such as FTP can generate a sustained file transfer process that can exist for tens of seconds and cause the traffic correlation to decay very slowly in time. The relative burst sizes generated by different applications is, therefore, an important consideration in the aggregation of such traffic.

To effectively implement differentiated and controlled load services at a network edge, it is necessary to identify the level of traffic aggregation that allows a robust traffic characterization that can be used for controlling the performance of queues. The approach presented here is to filter out from the aggregate traffic stream applications that exhibit significant nonstationary behavior and then fit the filtered stationary process with a parametric time series model. Such a model will allow for traffic prediction and forecasting. An example of a data traffic shaper that uses the time series parameters for controlling the correlation structure of the aggregate stream is demonstrated.

In Section 2.0, the relevant statistical features of the Internet traffic measurements analyzed in this work is presented. A stationarity test is used to identify applications that will require to be filtered from the aggregate traffic stream. Section 3.0 describes the fitting of a nonlinear time series model to the filtered traffic process. Section 3.1 compares the queuing performance of the simulated traffic with measurements. An application of the time series parameters in a traffic shaper is presented

in Section 4.0. Section 5.0 concludes the paper.

## 2. Statistical Features of Internet Traffic

The traffic measurements analyzed in this work are obtained at the junction where Ohio State University (OSU) connects to the vBNS (very high-performance Backbone Network Services). The data represent the aggregate traffic flowing from the OSU local area networks towards the vBNS backbone. The vBNS consists of an OC12 backbone and interconnects NSF-supported supercomputer centers, research and educational institutions and government sponsored networks. The measurements are collected using the OC3MON utility [4]. OC3MON is a programmable data collection and analysis tool developed by MCI and the National Laboratory for Networking Research (NLANR). Two network interface cards independently capture the traffic received and transmitted by the switch connecting OSU to vBNS. The OC3MON software extracts headers from packet traces and applies timestamps. The trace files are made available to the public at the archive *http://moat.nlanr.net/Traces*. These trace files are typically two minute duration samples. For this analysis an additional two hour trace was obtained from NLANR.

### 2.1. Long-term Traffic Trends

The traffic that flows in the direction from OSU to vBNS is analyzed. First the long-term statistics of two weeks of data from March 1, 1999 to March 14, 1999 are examined. This is motivated in part by the need for detection of invariant traffic features that can be used for high level traffic management. The daily traffic pattern is represented by eight measurement sets each representing approximately two minutes of traffic and obtained at three hour intervals. On examination of the total traffic volume generated in each measurement set it was evident that the peak usage typically occurs around 4 p.m and the traffic dips to a minima in the 4 - 7 a.m time frame. Other consistent features on the hourly time scale were the dominance of TCP traffic typically contributing over 95% of the total traffic volume throughout the day and followed by UDP protocol traffic at the 2-5% level which drops in volume significantly during the off-peak hours. Since TCP controls the queueing performance further analysis is focussed on the TCP packets.

The traffic measurements provide the time-stamp, packet size and protocol type, information on the source and destination ports and addresses. Further traffic classification is based on the source port numbers which represents the application being transported by TCP. The well known applications in these data files are Hypertext Transport Protocol (HTTP, port 80), File Transfer Protocol (FTP, port 20), Network News Transfer Protocol (NNTP, port 119) and Simple Mail Transfer Protocol (SMTP, port 25). Figure 1 displays the cumulative percentage of TCP traffic volume in bytes that is generated by each of these protocols around the 4p.m. peak hour in each day. The horizontal axis represents the measurement date. The well-known protocols typically contribute 70-80% of the total traffic volume during the week days. However, the combination of the less well-known *Other* applications are seen to also have a non-negligible influence in byte count to the aggregate traffic. These individual applications are generated by different source dynamics, packet sizes and interarrival time scales. In addition the traffic burst sizes may exhibit a wide variance between applications. An indiscriminate aggregation of these traffic streams may not necessarily lead to a homogeneous flow that optimizes the use of network bandwidth. To further characterize this traffic mix the application level traffic statistics are examined on a finer time scale by considering a continuous trace of two hour duration.
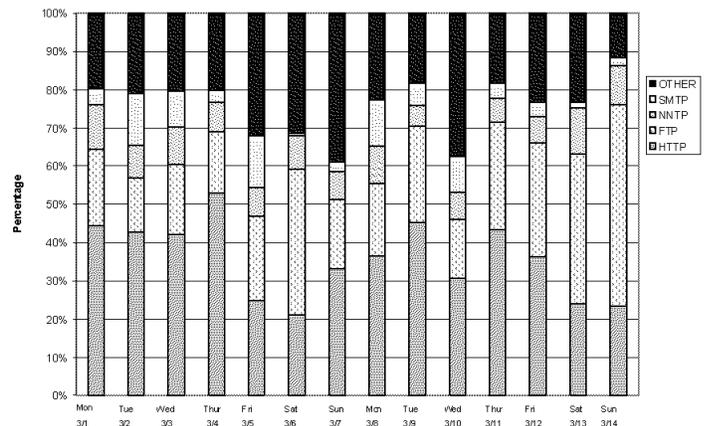


**Fig. 1** Traffic distribution per TCP application during peak usage (4 pm).(3/1/99 - 3/14/99)

### 2.2. Statistics of the Packet Arrival Process

To facilitate the analysis of individual application streams, the two hour packet trace was mapped to a time series by aggregating the packet arrivals over successive

non-overlapping time periods of duration 0.1 second. Traffic aggregation schemes must be designed to minimize queueing delays and losses while maintaining a reasonable level of network utilization. When stationary traffic arrival processes with randomly distributed burst patterns are multiplexed, it can be shown that the variance of the superposed stream reduces as the number of multiplexed sources is increased. This feature is desirable since it serves to improve the queueing performance in a switch buffer. If however, the arrival processes deviate from the properties of stationarity and randomness traffic aggregation may not lead to improved performance. In addition to stationary features the dependence in traffic stream is also important, since this impacts the approach taken towards traffic modeling. Typically traffic autocorrelations are examined for proposing models. This statistic fails to indicate presence of nonlinear features in the data which may also lead to traffic nonstationarity. Therefore the basic characteristics of stationarity and linearity of the time series is examined at the application level.

**2.2.1. Stationary Test** The traffic stationarity is determined by examining the variation of a sequence of arrival rates estimated using non-overlapping windows of fixed duration. The window size determines the time-scale of stationarity. For a fixed window size $\tau_s$, the set of mean arrival rates $\{a_i\}$, $i = 1, 2, \ldots N$ were determined, where $N$ represents the number of window segments that result from the time series. The sequence of $a_i$ for the aggregate time series may be examined visually for underlying trends and tested quantitatively under the hypothesis that each $a_i$ is an independent sample value of a random variable. Under this hypothesis the variations in the sample values $a_i$ will be random and exhibit no trends. Then the number of runs in the sequence relative to, say the median value will be as expected for a sequence of independent observations of the random variable [5]. If the probability of remaining below or above the median value remain unchanged from one $a_i$ to the next, then the sampling distribution of the number of runs in the sequence of $a_i$ may be determined and tabulated for various confidence intervals. This is referred to as the run-test. In applying this test to the traffic data, the temporal correlation of the time series must be taken into consideration and the selection of the time window $\tau_s$ should be governed by the maximum correlation time scale we wish to accomodate in our traffic model. The value of $\tau_s$ is selected in the range of $36 - 72$ seconds, yielding 200 and 100 values of $a_i$ respectively. For each $\tau_s$, the number of runs of $a_i$ rela-

tive to the median were determined and evaluated against the expected number for a 95% confidence interval.

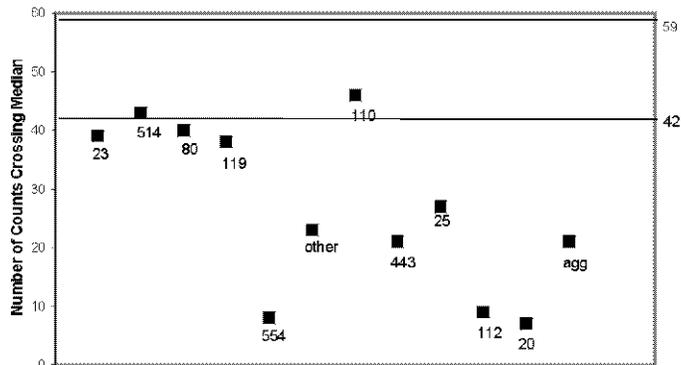The results of the run test for $N = 100$ are shown in Figure 2.



**Fig. 2** Stationarity test(numbers are application port ids; agg:aggregate traffic)

The two solid lines depict the range from [42-59] for which the stationarity hypothesis would be accepted. The stationarity test was conducted for the aggregate time series denoted *AGG* and for each individual time series constructed for the applications. These applications are represented in Figure 2 using the port numbers. The application ports highlighted in the figure are the top ten applications that contribute to over 85% of the TCP traffic process. It is seen that while the aggregate TCP process exhibits significant deviation from stationarity, the applications corresponding to the ports { 80, 119, 23, 514, 110 } approach the boundaries of being accepted as stationary processes. Applications such as FTP (20) exhibit significant deviation from stationarity. A closer examination of the FTP time series showed the presence of one long file-transfer burst existing for over ten seconds, leading to a strong contribution to the total traffic and deviating from random process properties. The breakup in traffic percentage of each application is tabulated in Table I. For the remainder of this paper, we will consider the analysis and modeling of the stationary process obtained by aggregating the applications that individually satisfy conditions of stationarity.

**Table 1:** Percentage of traffic from different applications

| Port ID | Traffic Fraction (%) |
|---------|----------------------|
| 80 | 43.37 |
| 20 | 19.87 |
| 119 | 10.84 |
| 25 | 6.11 |
| 112 | 4.44 |
| 110 | 0.91 |
| 443 | 0.75 |
| 554 | 0.56 |
| 23 | 0.49 |
| 514 | 0.29 |

**2.2.2. Linearity Test** We consider the TCP traffic arrival process obtained by filtering out application processes identified to exhibit nonstationary features. The normalized autocorrelation function (NACF) of this filtered process may be compared to that of the aggregate
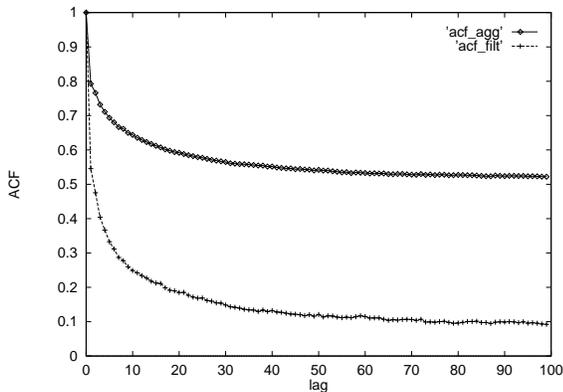


**Fig. 3** Comparison of ACFs of aggregate and filtered traffic processes.

TCP stream in Figure 3. The removal of nonstationary components is seen to lead to a significant reduction in the temporal correlation of the packet arrival process. The NACF is a measure of linear dependence between the random variates of the time series and can suggest a choice of linear predictive models such as autoregressive and/or moving average (ARMA) models that can capture the observed correlation function. These models generally require that the model errors be Gaussian distributed. Therefore tests of linearity and normality of the data will have to be confirmed before such models are hypothesized.

To verify linearity and in the process rule out nonlinear dependence conditional statistics as given by sample regression functions [6] are obtained. The lag $j$ re-

gression function of a time series $X_n$ is defined as $r_j = E[X_n | X_{n-j} = x]$. The estimates of $r_j$ for $j = 1, 2, 3, 4$ are depicted in Figure 4. The horizontal axis represents a partition of the amplitude range of the time series into a finite set of disjoint sets. The vertical axis represents the function $r_j$ determined by calculating an average of all samples that satisfy the regression constraint. For Gaussian distributed linear processes, the regression functions exhibit a linear trend. The regression functions depicted in Figure 4 show a deviation from linear structure. It was found that if one ignores the nonlinear dependence and fits the best linear ARMA model to the measurements, this model fails to capture the variability and the index of dispersion of the arrival process and significantly underestimates the bit loss statistics in a queue. Figure 4 indicates that a piece-wise linear model may be more appropriate to capture the dependence structure in the time series amplitudes. In particular, the regression functions suggest that the amplitudes below and above a chosen threshold level exhibit differing correlation structures as evidenced by the change in their slopes. This is important from a network operating perspective since it indicates that traffic features in the state leading to congestion may be differentiated from that of more benign states, allowing the design of more robust predictive models. In the next section, we propose a nonlinear time series model that can describe the observed traffic nonlinearities.
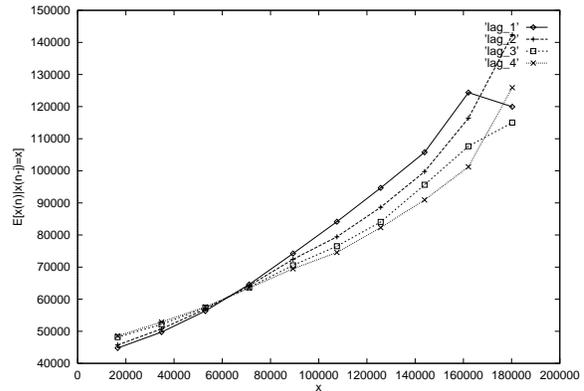


**Fig. 4** Regression functions $r_j$  $j = 1, 2, 3, 4$

## 3. Threshold Autoregressive (TAR) Models

The threshold autoregressive model (TAR) is proposed and is a nonlinear model comprised of a set of linear AR models which are valid in disjoint subregions of the time series amplitude. At a given time the esti-

mated traffic amplitude will be based on the AR process that governs a particular subregion. This subregion is selected based on the amplitudes observed over previous time values. These lagged observations serve as a decision criteria for generating the current traffic estimate. The TAR model and its variants have been successfully applied for modeling time series data exhibiting cyclical properties and long-range dependence [7].

The model considered here will incorporate two amplitude ranges denoted as low ($L$) and high ($H$) amplitude states. In the low state the traffic takes on values $L: (0, \hat{r}]$, where $\hat{r}$ is a threshold value. The high state accommodates amplitudes $H: [\hat{r}, \infty)$. In addition a delay value $d$ will be used to determine the state of the process. The combination of two amplitude regimes and a single delay will result in two subregions that describe the time series. These subregions are denoted as $R_j$ $j = 1, 2$. The governing AR process $x(n)$ at time $n$ is selected based on the amplitude of the time series at the lagged amplitude $x(n - d)$.

For each of the $R_j$ the process evolves as a stable AR process, governed by the correlations within that region. The current value of the byte rate at time $n$ will be governed by an autoregressive process of order $k_j$.

$$x(n) = \alpha_0^j + \sum_{i=1}^{k_j} \alpha_i^j x(n - i) + e^j(n) \qquad (1)$$

when $x(n - d) \in R_j$. Here the $\alpha_0^j = \mu^j \left[1 - \sum_{i=1}^{k_j} \alpha_i^j\right]$ where $\mu^j$ and $\alpha_i^j$, $i \geq 1$ represent mean value and AR coefficients of the stationary process in subregion $R_j$. The term $e^j(n)$ represents residual errors assumed to be derived from an independent identically distributed random process having a zero mean and finite variance $\sigma^{j^2}$. When subregion constraints are violated, the process is switched to the subregion model that obeys the proper amplitude and delay constraints. The delay parameter affords the flexibility of capturing persistence phenomenon at the required amplitudes. Extended sojourn in the high byte rate state is an important feature in delay management, whereas dwell-time in the low byte-rate state has impact on multiplexing efficiency. Therefore, the thresholds and delay parameters should be carefully chosen to capture these critical elements in the traffic variation.

## 3.1. TAR Model Parameter Selection

To construct the model, the optimal values of $\hat{r}$

and the delay $d$ must be selected along with the coefficients of the local AR processes for each subregion. This is accomplished by varying $d$ and $\hat{r}$ over a selected range and for each such pair, the local AR coefficient vector $\underline{\alpha}^j$, and the residual error variance $\sigma_j^2$ are determined using least squares estimators. The process of TAR parameter estimation is that outlined in Tong [8].

To estimate the local AR parameters, the data is searched for all samples $x(n)$ that satisfy the given amplitude and delay constraints $R_j$. For each $R_j$, the vector $\underline{x}^j = (x_{i_1}^j, \ldots, x_{i_{n_j}}^j)^T$ contains the $n_j$ time series samples that satisfy constraint $R_j$. For this data, the $k_j^{th}$ order linear model coefficients are evaluated.

$$\underline{x}^j = A^j \underline{\alpha}^j + \underline{e}^j \qquad (2)$$

where $\underline{\alpha}^j: (\alpha_0^j, \alpha_1^j, \ldots, \alpha_{k_j}^j)^T$ and $A^j$ is a $n_j \times (k_j + 1)$ matrix with each row comprised of the $k_j$ values that lag the elements of $\underline{x}^j$ and $\underline{e}^j = (e_{i_1}^j, \ldots, e_{i_{n_j}}^j)^T$ is the residual error vector. Since $n_j \gg k_j$, the solution of the system of equations in Eq. (2) is the least squares estimate of $\underline{a}^j$, denoted as $\underline{\hat{a}}^j$. The error variance $\sigma_j^2 = ||\underline{\hat{e}}^j||^2 / n_j$ is determined as the approximate Maximum Likelihood Estimate of the noise variance.

For the fixed delay $d$, the threshold amplitude $\hat{r}$ is sampled uniformly in the range 50000 and 70000 bytes in intervals of 5000 bytes. This range is selected to vary about the average value of the time series. The delay was varied from one to four lags. The maximum allowable model order of the AR process was selected to be 30, allowing for correlations upto three second time scales. For each subregion $R_j$ the least squares estimates of the local AR coefficients was determined using Eq. (2). The optimum order $k_j$ for the $j^{th}$ submodel corresponds to the value $k$ that yields the minimum value for the Akaike Information Criteria (AIC) statistic $AIC(k_j)$ [9]. This process is repeated for all subregions, for each value of $\hat{r}$ and $d$. The total AIC as a function of the threshold and delay parameters is computed as $AIC_{total}(d, \hat{r}) = \sum_{j=1}^{2} AIC(k_j)$. The optimal threshold parameter $\hat{r}$, delay $d$ and AR parameters are those that yield the minimum $AIC_{total}(d, \hat{r})$. This process resulted in $d = 1$, $\hat{r} = 70000$, $k_1 = 21$, $k_2 = 9$. Although these appear to be high order processes, an examination of the magnitudes of $\alpha_i$ reveals that only the first three AR coefficients have dominant magnitudes.

To simulate the TAR process, the initial condi-

tions were determined from the measurements. Subsequent values were generated using the TAR parameter estimates and additive noise variables for each subregion. The noise process was derived from the empirical distributions of the error residuals obtained during the fitting of the measurements to the model. These distributions were tested for Gaussian statistics and found to pass this test at acceptable significance levels. In the simulation, any negative values that resulted were set to be equal to a preset minimum as determined from the measurements. Results from the TAR model simulation are presented next.

### 3.2. Model Evaluation

Several statistical features were derived from the simulated data and measurements to evaluate the performance of the traffic model. The QQ (quantile-quantile) plots of marginal distributions of the byte-amplitudes and autocorrelation functions are shown in Figures (5) and (6). The visual assessment of the QQ plots indicate a good agreement in the marginal distributions. The acf plots show that the TAR model captures the trend in the data up to about 50 lags, a five second timescale. The TAR model was also found to perform well in capturing the counting statistics as evaluated by the expected value and the variance in the number of arrivals as a function of time up to timescales of five seconds.

To determine if an adequate number of timescales were represented in the model, the performance of the model was evaluated in a finite buffer queue, serviced at a constant rate. The results of the bit-loss-ratio (BLR) are plotted in Figure 7. The horizontal axis represents the service rate normalized by the average rate. The bit loss ratios are depicted for buffer sizes varying from one millisecond to fifty milliseconds. These results are compared to the losses experienced by driving the queue with the measurements. The results show that the TAR model exhibits reasonable agreement in the loss statistics in this range of buffer sizes.

### 4. TAR Model Based Traffic Shaper

The TAR model has been shown to provide a simple and reasonably accurate model for simulating data traffic patterns of the dominant TCP applications such as http. One of the goals in this traffic analysis work is to derive implementable algorithms for controlling the performance of queues that integrate real-time (RT) delay sensitive and delay tolerant traffic. In this context, TCP data traffic may be considered to be delay tolerant. One approach for controlling the performance

is to shape the TCP traffic to reduce its impact on the RT flows. A traffic shaper that is parametrized by the TAR traffic model parameters is proposed. The traffic shaper is designed to reduce the impact of TCP traffic correlations on the RT traffic. It is well known that positive correlations in the arrival process cause performance degradation [10]. The time series model provides a parametric representation of the correlation structure which allows the design of filters that can decorrelate the arrival process, prior to the admission into a shared buffer. The process of removing all correlations in the data leading to a sequence of independent random variates is also referred to as a whitening filter.

The traffic shaping application is depicted in Fig.8(a), where $x(n)$ represents the measurement trace aggregated into a time series at time $n$ and $c(n)$ is the time-varying drain rate of the shaping buffer. The filter is an $m^{th}$ order linear predictor of $x(n)$, obtained from the lagged input samples as, $\hat{x}(n) = \sum_{i=1}^{m} \alpha_i^j \, x(n-i)$. Here $m \leq k^j$ the AR order of the input signal in subregion $R_j$. The subregion identification for $x(n)$ is obtained based on the amplitude of $x(n-1)$. The admission rate $C(n) = \mu + [x(n) - \hat{x}(n)]$ is the rate at which TCP traffic is admitted into the common buffer at time $n$ and varies about a fixed mean arrival rate $\mu$. If $m < k^j$, the output process $C(n)$ is characterized by a reduced order correlation relative to that of $x(n)$. The higher the value of $m$ the larger will be the delays experienced by the data in the shaping buffer. The correlation of the shaped traffic stream for $m = 3$ and that of the measurement trace $x(n)$ is depicted in Fig.8(b). By shaping the bursts of the data traffic process in this manner, the multiplexing of data with delay sensitive traffic in a common buffer may be accomplished while controlling the quality of service constraints.

### 5. Conclusions

Statistical analysis of data traffic measurements originating from the local area networks of a campus site has identified that a subset of TCP based applications induce nonstationary features in the aggregate traffic stream. This feature is problematic for the prediction and forecasting of traffic statistics. These applications which consist of FTP, SMTP and other less known ports are identified through stationarity tests and filtered from the aggregate process. The resulting traffic stream comprises about 60-70% of the traffic and consists primarily of the well-known http and nntp applications. This filtered traffic stream is modeled as a nonlinear threshold

autoregressive process. The proposed model captures the correlation structure in the measurements and shows reasonable agreement with the packet loss statistics. The parametric model is shown to be applicable in the design of traffic shapers that can decorrelate the data traffic and allow delay management in queues shared with delay sensitive traffic.

## Acknowledgements

**Fig. 7** BLR comparison between simulation (symbols) and data. Buffersize= 0.001, 0.01, 0.05 (secs).



**Fig. 5** QQ plot comparison of TAR model and data.



**Fig.8(a)** Linear prediction based traffic shaper.



**Fig. 6** Comparison of ACFs of TAR model and data.
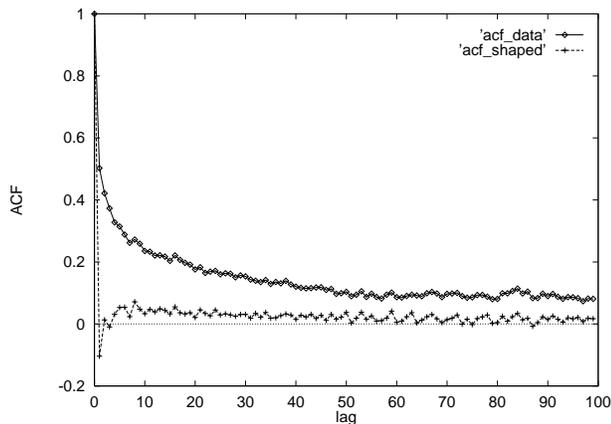
**Fig.8(b)** Comparison of NACFs of data and shaped traffic.

## References

1. W.E. Leland, M.S. Taqqu, W. Willinger and D.V. Wilson, 'On the Self-Similar Nature of Ethernet Traffic (Extended Version),' IEEE/ACM Trans. Networking, p1-15, **2** (1), 1994.

2. K. Meier-Hellstern, P.E. Wirth, Y.-L. Yan and D.A. Hoeflin, "Traffic Models for ISDN data users: Office Automation application," in *Teletraffic and Data Traffic in a Period of Change,* Eds. A. Jensen and V.B. Iversen, Proc. of ITC-13, Copenhagen, p167-172, Elsevier Science Pub., Amsterdam, 1991.

3. V. Paxson and S. Floyd, "Wide-Area Traffic: The Failure of Poisson Modeling", IEEE/ACM Trans. Networking, **3**, p226-244, June 1995.

4. J. Apisdorf, K. Claffy, K.Thompson and R. Wilder,"OC3MON: Flexible, Affordable, High Performance Statistics Collection," $10^{th}$ Systems Administration Conference (LISA'96), USENIX, Sept.29-Oct. 4, Chicago, IL., 1996.

5. J.S. Bendat and A.G. Piersol, **Random Data: Analysis and Measurement Procedures,** Chapt. 7, Wiley-Interscience, 1971.

6. H. Tong, **Threshold Models in Nonlinear Time Series Analysis,** Lecture Notes in Statistics, vol. 21, Springer-Verlag, 1983.

7. P.A.W. Lewis and B.K. Ray, "Modeling Long-Range Dependence, Nonlinearity and Periodic Phenomenon in Sea Surface Temperatures using TSMARS," J. American Statistical Association, **92** (439), p881-893, September 1997.

8. H. Tong, **Nonlinear Time Series , A Dynamical System Approach,** Oxford Science Publications, Clarendon Press, Oxford, 1990.

9. S.M. Kay, **Modern Spectral Estimation,** Chap. 7, Prentice Hall, Englewood Cliffs, NJ, 1988.

10. M. Livny, B. Melamed and A.K. Tsiolis, "The impact of autocorrelation on queuing systems," *Management Science,* 39, p322-339, (1993).