A Neural-Based Technique for Estimating Self-Similar Traffic Average Queuing Delay

Homayoun Yousefi'zadeh, Member IEEE

Abstract—Estimating buffer latency is one of the most important challenges in the analysis and design of traffic control algorithms. In this paper a novel approach for estimating average queuing delay in multiple source queuing systems is introduced. The approach relies on the modeling power of neural networks in predicting self-similar traffic patterns in order to determine the arrival rate and the packet latency of low loss, moderately loaded queuing systems accommodating such traffic patterns.

Index Terms—Bursty Traffic, Self-Similarity, Intelligent Traffic Modeling, Neural Network, Queuing Delay.

I. INTRODUCTION

NALYSIS of traffic data from networks and services such as Ethernet LANs [6], Variable Bit Rate (VBR) video [1], ISDN traffic [4], and Common Channel Signaling Network (CCNS) [2] have all convincingly demonstrated the presence of features such as long range dependence, slowly decaying variances, and heavy-tailed distributions. These features are best described within the context of second-order self-similarity and fractal theory. Self-similar phenomena show structural similarities across a wide range of time scales in which traffic spikes ride on the longer term ripples, that in turn ride on longer term swells, so on and so forth.

Neural networks are a class of nonlinear systems capable of learning and performing tasks accomplished by other systems. Some of the applications of neural networks manifest in speech and signal processing, pattern recognition, and system modeling. Systems with neural network building blocks are robust in the sense that occurrence of small errors does not interfere with proper operation of the system. This characteristic of the neural networks makes them quite suitable for traffic modeling.

Estimating and reducing packet latency is a major design issue in computer communication networks. There are a number of factors that introduce delay in network services. Different delay types may be classified under processing, propagation, multiplexing, and queuing categories. The main objective of packet scheduling methods is then to come up with solutions for predicting and reducing delay while efficiently utilizing network resources.

In [8], we made use of the modeling power of neural networks introduced in [7] to provide a fair dynamic buffer management scheme improving the loss performance of a class of queuing systems with self-similar characteristics. In this study, we utilize the modeling power of neural networks in predicting self-similar traffic patterns in order to determine the arrival rate and the packet latency of queuing systems accommodating such patterns assuming the service rates are given and there is no significant loss impact. Our packet latency estimation technique might be thought of as a part of a packet scheduling algorithm.

An outline of the paper follows. Section II briefly reviews the characteristics of aggregated self-similar traffic patterns and provides an overview of the neural network modeling of such traffic patterns. Section III describes typical multiple source systems used for the application task and discusses the packet latency estimation application. It also evaluates the performance of the average latency estimation technique by comparing its results with measured average latency in the presence of typical buffer management and server scheduling schemes. The paper concludes in Section IV.

II. SELF-SIMILAR TRAFFIC MODELING

This section includes a brief description of self-similarity followed by a review of the neural network modeling technique.

A. Second-Order Self-Similarity

In [7] and [8], we provide an analytical framework for selfsimilarity as a statistical property of the time series. Mathematically, self-similarity manifests itself in a number of ways.

- Slowly decaying variance property points out that the variance of sample mean decreases more slowly than the reciprocal of the sample size with the meaning $var(X^{(m)}) \sim k_2 m^{(-\beta)}$ as $m \to \infty$ for $0 < \beta < 1$.
- Long range dependence property expresses that the autocorrelations decay hyperbolically rather than exponentially fast, implying a non-summable autocorrelation function, i.e., $\sum_{n} R(n) = \infty$.
- 1/f noise property mentions that the spectral density f(.) obeys a power-law near the origin with the meaning f(λ) = k₃λ^{-γ} as λ → ∞ for 0 < γ < 1 and γ = 1 − β.

The most important feature of self-similar processes is seemingly the fact that their aggregated processes $X^{(m)}$ possess a non-degenerate correlation function as $m \to \infty$. This is completely different from typical packet traffic models previously considered in the literature all of which have the property that their aggregated processes $X^{(m)}$ tend to second order pure noise, i.e., $R^{(m)} \to 0$ as $m \to \infty$.

B. Neural Network Modeling of Self-Similar Traffic

In [7], we describe how a fixed structure feed forward perceptron neural network with back propagation learning algorithm can be used to model aggregated self-similar traffic patterns as

The author is with the Center for Pervasive Communications at the Electrical and Computer Engineering Department of University of California, Irvine; (e-mail:hyousefi@uci.edu).

an alternative to stochastic and chaotic systems approaches proposed in [5] and [3]. We note that although the emphasis of our work is on self-similar traffic modeling, our proposed neural network modeling approach can nevertheless be used for any traffic pattern independent of self-similarity. In what follows we briefly review the neural network modeling technique of [7] in which an elegant approach capable of coping with the fractal properties of the aggregated traffic is introduced. The approach provides an attractive solution for traffic modeling and has the advantage of simplicity compare to the previously proposed approaches namely stochastic and deterministic chaotic map modeling. The promise of neural network modeling approach is to replace the analytical difficulties encountered in the other modeling approaches with a straight forward computational algorithm. As oppose to the other modeling approaches, neural network modeling does not introduce a parameter describing the fractal nature of traffic neither does it investigate identification of appropriate maps. It, hence, need not cope with the complexity of estimating Hurst parameter and/or fractal dimensions. The approach simply takes advantage of using a fixed structure nonlinear system with a well defined analytical model that is able to predict a traffic pattern after learning the pattern dynamics through the use of information available in a number of traffic samples.

The fixed structure, fully connected, feed forward perceptron neural network utilized for the task of modeling consists of an input layer with eight neurons, three hidden layers with twenty neurons in each layer, and an output layer with one neuron. Figure (1) illustrates the structure of the neural network. The output of each neuron is connected to the input of all of the neurons in the layer above after being multiplied in a weighting function. The specific neural network used for the task of modeling relies on the so-called back propagation learning algorithm described in [7] and the references therein. In a nutshell, the back propagation learning algorithm changes the weighting functions of the underlying neural network in the opposite direction of the gradient vector and its momentums in order to minimize the absolute error function defined proportional to the square of the difference between the neural network output and the real output.

In a typical iteration of the learning phase, the neural network is provided with samples x[k-8] through x[k-1] of the real traffic pattern and the difference between sample x[k] of the real traffic pattern and the neural network output is used to adjust the weighting functions of the network accordingly. In the next iteration, sample x[k-8] of the real traffic pattern is discarded, samples x[k-7] through x[k] of the real traffic pattern are used as the new input sample set, and sample x[k+1] is used as the new real traffic sample. The neural network continues processing more information in consecutive iterations of the learning phase until the absolute error is less than a specified error bound, ϵ . The learning phase of the perceptron neural network is directly followed by the recalling phase when the network output is able to follow the real traffic within the acceptable error bound, ϵ . In each iteration of the recalling phase, the neural network independently generates the samples by discarding the oldest input sample, shifting the input samples by one, and using its output as the most recent input sample. The same se-

quence of following a learning phase by a recalling phase is repeated when and if the neural network output difference exceeds the acceptable error bound, ϵ . The number of samples required for the training of the neural network depends on the complexity of the traffic pattern dynamics. The time complexity and the space complexity of the back propagation algorithm are respectively $\mathcal{O}(IN)$ and $\mathcal{O}(N)$ where N is the number of weighting functions in the network and I is the number of iterations. Although the complexity is typically better than the complexity of implementing statistical approaches such as fractional ARIMA processes or the complexity of calculating fractal dimensions such as correlation dimension, wide variations of I prevent us from making a strong claim about complexity advantage of the algorithm compare to other algorithms. Nonetheless combining the straight forward way of implementation with the analysis of complexity, we claim that the neural network modeling approach provides an elegant approach for the task of traffic modeling.



Fig. 1. Fixed structure neural network used for the task of modeling.

In the following section, we apply the proposed neural network modeling technique to predict the packet generation patterns of a number of ON-OFF traffic sources and utilize the prediction results in estimating arrival rates and average latencies in queuing systems accommodating such patterns.

III. LATENCY IN SELF-SIMILAR QUEUING SYSTEMS

Our application test bed relies on a multiple source queuing system. A multiple source queuing system consists of a number of sources sharing a total available buffer space. In [8], we provide a brief queuing analysis for individual queues of such a multiple source queuing system. Traffic pattern of each source includes the packets generated by a number of ON-OFF chaotic maps. An ON-OFF source model is generating traffic at a peak rate when it is active and becomes active as soon as the state variable of the describing chaotic map goes beyond a threshold value. The source becomes passive as soon as the state variable goes below the threshold value. We utilize double intermittency map in our packet generation process as it generates a self-similar traffic pattern according to what is described in [3]. We propose using different initial conditions for a fixed threshold value to obtain different traffic patterns for different sources. As an alternative, one may use different threshold values with fixed or variable initial conditions to achieve varying traffic patterns for different sources.

We now apply our neural network modeling scheme to predict the total number of generated packets and utilize the prediction results in estimating the queuing delay for the packets generated by a number of traffic sources. Consider a multi-



Fig. 2. The structure of a multiple source queuing system.

ple source queuing system such as the one shown in Figure (2) with three sources sharing the space available in a central buffer. Assume that the aggregated traffic pattern of each individual source consists of the traffic patterns of 120 sources generating ON-OFF packet traffic according to double intermittency map model. The buffer size is assumed to be fixed, large enough to prevent any loss. In addition, suppose that the system is utilizing complete sharing buffer management, Statistical Time Division Multiplexing (STDM) scheduling, and First Come First Serve (FCFS) service discipline schemes as described in [8]. Utilizing the prediction results of our neural network modeling scheme, we can estimate the packet arrival rate of the central buffer. For a given service rate and a known buffer occupancy, the queuing delay of a packet can be measured as the average number of time units it spends in the queue before leaving the buffer.

Figure (3) displays our simulation results for the system described above. It shows the Measured Average Latency (MAL) and the Estimated Average Latency (EAL) versus service time diagram for the triple source queuing system over the intervals in which the arrival rate predictions are of acceptable accuracy. The average latency has been calculated over the time periods in which the neural network has been able to follow the arrival pattern of the central buffer. For the relative error defined as $\frac{|MAL-EAL|}{MAL}$, Figure (3) shows that the estimation results are MALwithin the 3% relative error range pending the following conditions are held. First, the averaging period is long enough in order for the neural network to be able to follow the traffic pattern for a number of times within the specified error bounds and second, the buffer service rate does not exceed an existing threshold value. Although not shown in the simulation results, we have observed that the average packet latency drops sharply by choosing service rates beyond the threshold value. In the latter case, the neural network latency estimation findings are not acceptable as the result of having high service rates and low average latencies. The service rate threshold generally depends on the dynamics of the system and for the triple source system of our experiment is the normalized value 13.

We finish this section by mentioning that a typical sequence of learning and recalling phases consists of few hundred thousand samples and hundreds of samples respectively. In addition, all of the convergence results are strongly affected by the choice of initial conditions of the weighting functions of the



Fig. 3. Estimated average latency (EAL), and measured average latency (MAL) versus normalized service rate for the triple source queuing system.

neural network. As a practical finding, setting the initial values of the weighting functions of the neural network at 0.01 typically yields good results. Our justification for both of the above phenomena is the fact that the proposed neural network is trying to learn complicated dynamics of chaotic maps exhibiting extreme sensitivity to variations of initial conditions.

IV. CONCLUSION

In this paper, we introduced a novel approach for estimating queuing latency in multiple source queuing systems as an application of neural network modeling of self-similar packet traffic. We relied on the prediction power of neural networks to estimate arrival rates and packet latencies in multiple source queuing systems accommodating self-similar traffic patterns. We evaluated the performance of our estimation technique by comparing estimated average latency with measured average latency and concluded that the scheme is able to provide an acceptable estimate with a less than 3% relative error below a specified service rate threshold for moderately and heavily loaded systems with no significant loss.

REFERENCES

- J. Beran, R. Sherman, M. S. Taqqu, W. Willinger, "Variable Bit Rate Video Traffic and Long Range Dependence", IEEE/ACM Trans. on Networking, Vol. 2, NO. 3, Apr. 1994.
- [2] D. E. Duffy, W. Willinger, "Statistical Analysis of CCSN/SS7 Traffic Data from Working CCS Subnetworks", IEEE JSAC, 1994.
- [3] A. Erramilli, R. P. Singh, P. Pruthi, "Chaotic Maps as Models of Packet Traffic.", ITC Vol. 14, PP. 329-338, 1994.
- [4] K. M. Hellstern, P. Wirth, "Traffic Models for ISDN Data Users: Office Automation Application", Proc. ITC-13, Denmark, 1991.
- [5] W. E. Leland, W. Willinger, M. S. Taqqu, D. V. Willson, "Statistical Analysis and Stochastic Modeling of Self-Similar Datatraffic", ITC Vol. 14, PP. 319-328, 1994.
- [6] W. E. Leland, W.Willinger, M. S. Taqqu, D. V. Willson, "On the Self-Similar Nature of Ethernet Traffic", IEEE/ACM Trans. on Networking, Vol. 2, NO. 1, PP. 1-15, Feb. 1994.
- [7] H. Yousefi'zadeh, "Neural Network Modeling of a Class of ON-OFF Source Models with Self-Similar Characteristics", Submitted to IEEE/ACM Trans. on Networking, February 2002. Also available at http://www.ece.uci.edu/~ hyousefi/pub.html.
- [8] H. Yousefi'zadeh, E. A. Jonckheere, J. A. Silvester, "A Fair Dynamic Neural-Based Buffer Management Scheme for Improving Packet Loss Performance in a Class of Queuing Systems with Self-Similar Characteristics", Submitted to IEEE/ACM Trans. on Networking, March 2002. Also available at http://www.ece.uci.edu/~ hyousefi/pub.html.