

An application of neural networks to adaptive playout delay in VoIP

Ying Zhang¹, Damien Fay¹, Liam Kilmartin²

¹Department of Mathematics, National University of Ireland, Galway
damien.fay@nuigalway.ie

²Department of Electronic Engineering, National University of Ireland, Galway
liam.kilmartin@nuigalway.ie

Abstract

The statistical nature of data traffic and the dynamic routing techniques employed in IP networks results in a varying network delay (jitter) experienced by the individual IP packets which form a VoIP flow. As a result voice packets generated at successive and periodic intervals at a source will typically be buffered at the receiver prior to playback in order to smooth out the jitter. However, the additional delay introduced by the playout buffer degrades the quality of service. Thus, the ability to forecast the jitter is an integral part of selecting an appropriate buffer size. This paper compares several neural network based models for adaptive playout buffer selection and in particular a novel combined wavelet transform/neural network approach is proposed. The effectiveness of these algorithms is evaluated using recorded VoIP traces by comparing the buffering delay and the packet loss ratios for each technique. In addition, an output speech signal is reconstructed based on the packet loss information for each algorithm and the perceptual quality of the speech is then estimated using the PESQ MOS algorithm. Simulation results indicate that proposed Haar-Wavelets-Packet MLP and Statistical-Model MLP adaptive scheduling schemes offer superior performance.

Keywords: VoIP, playout delay, neural networks, time series forecasting.

1. Introduction

In recent years Voice over IP (VoIP) has seen a huge increase in use due to its cost effectiveness, support of multimedia technology and ease of use. However, the network delay and packet loss, which are ubiquitous due to the best-effort mechanism on which significant portions of the internet are still based are the main factors affecting the Quality of Service (QoS) of a VoIP call[1]. When audio packets are transmitted over the internet, the variable network delay (jitter), which is mainly due to the variable queuing time in routers, modifies the periodic form of the transmitted audio packets when these packets are observed at the receiver [1] as is shown in Figure 1. The playout delay process is an application which aims to reduce the impact of network delay variability by buffering the received packets and playing them out after a certain time. Any packets which arrive later than their playout delay time are regarded as 'lost packets' and hence are not played out. Increasing the playout delay can reduce the packet loss, but a long playout delay has a negative impact on the real-time communication quality. Thus, a trade-off exists between the playout delay and packet loss rate. For interactive audio, a packet delay (due to all contributors of delay) of up to 400ms [2] and packet loss rate less than 5% are considered adequate [3].

In early VoIP system, a fixed playout delay was proposed as an initial solution to this problem [4]. While this method offers an easily implemented solution, it is not an optimum solution as it does not

take into account the fact that network jitter varies with time, as illustrated in Figure 1. Modern VoIP systems utilise adaptive playout delay approaches which estimate the network jitter continuously and dynamically adjusts the playout delay at the beginning of each talkspurt. Many algorithms have been proposed for estimating the network jitter such as Autoregressive (AR) models [5], Moving Average (MA) models [6], other statistical models [7-10], and adaptive filter models [11, 12]. In this paper, two new approaches based on combining Artificial Neural Network (ANN) and wavelet techniques and Artificial Neural Network (ANN) and Statistical Models are presented.

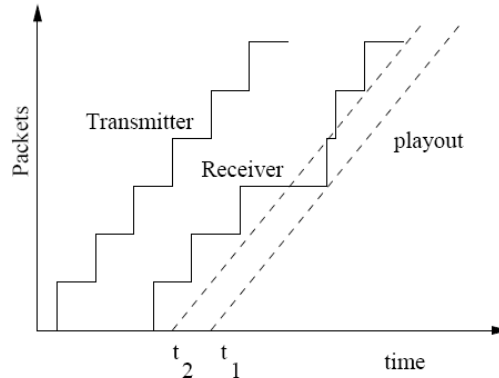


Figure 1 : Voice packets over network [4]

2. Proposed Models

The initial phase of this research focused on the use of neural networks as methodologies for the prediction of network delay jitter. Three types of neural network were compared for forecasting the network delay jitter in the network. The first two types are based on a standard multi-layer perceptron (MLP) and a recurrent-MLP [11] with standard training and cross validation methodologies been used to optimise network structure and performance. The last network utilises a wavelet transform as an input stage prior to applying the resultant signal to an MLP thus forming a Wavelet Packet-MLP (WP-MLP) [12]. Later research being presented focuses on utilising neural networks to predict the parameters of statistical models of the jitter waveform as an alternative approach.

The traditional back-propagation algorithm using Levenberg-Marquadt with cross validation has been used to train the networks [11]. The data is split into three different data sets used for training, validation (used for cross-validation and structure determination) and testing (used to compare each model *after* training). Each of the networks has two hidden layers and two outputs. To determine a suitable structure for the network (i.e. the number of nodes in each layer), different network structures were trained (ranging from a 2×2 to a 13×13 network) and their Prediction Mean Squared Errors (PMSE) compared over the validation set. The best structure was then selected for further evaluation. According to the PMSE performance on the validation set as shown in Table 1 below, the best performing structure was a 10×3 network with MSE 5.6×10^{-6} .

MLP Structure	Second Layer Nodes											
	2	3	4	5	6	7	8	9	10	11	12	13
First Layer	2	19.65										
	3	2.42	5.42									
	4	6.85	13.55	9.79								
Nodes	5	2.72	1.34	1.31	9.83							
	6	1.64	6.88	1.49	9.88	9.49						
	7	9.55	9.90	0.91	21.68	7.08	7.37					
Nodes	8	9.56	8.94	2.97	3.28	8.99	1.26	9.24				
	9	2.80	8.24	1.51	0.93	1.32	9.89	9.74	1.47			
	10	1.93	0.56	5.87	9.84	8.32	4.56	3.76	4.03	1.19		
	11	2.82	1.50	0.84	9.91	4.11	8.47	9.99	8.16	9.79	8.98	
	12	9.70	9.89	1.07	2.81	9.28	9.78	5.79	9.71	7.75	2.20	9.85
	13	4.31	0.54	0.36	1.27	9.49	2.31	7.08	9.26	4.68	8.64	7.04

Table 1 : PMSE of Different MLP Structures ($\times 10^{-5}$)

2.1 Wavelet-Packet Neural Network

In recent years, wavelet networks for function approximation and more specifically time series forecasting has been proposed in [12] and [13]. The wavelet transform maps a time domain signal into a time-frequency domain signal in which the coefficients represent the signal at progressively smaller frequency bands covering larger time spans. Specifically, given a discrete time series $x(k)$ the wavelet transform projects this series onto a new domain known as a *wavelet basis* [14], as:

$$x(k) = \sum_{i=0}^{+\infty} \sum_{j=-\infty}^{+\infty} w_{i,j}^y(k) \psi_{i,j}(k) \quad (1)$$

where $w_{i,j}^y(k)$ are called the *wavelet coefficients* defined as the inner product of $x(k)$ and the basis vectors $\psi_{i,j}(k)$:

$$w_{i,j}^y(k) = \int_{-\infty}^{+\infty} y(t) \psi_{i,j}^*(k) dt \quad (2)$$

and

$$\psi_{i,j}(k) = a_0^{i/2} \psi(a_0^i t - k \tau_0) \quad (3)$$

where $\psi(t)$ is called the *mother wavelet* and $\psi_{i,j}(k)$ is defined in terms of dilations (expansion), a_0 , and translations (phase shift), τ_0 , of the mother wavelet and * denotes the complex conjugate. There are various types of mother wavelet, such as Haar wavelet, Meyer Wavelet, Coiflet wavelet, Daubechies wavelet, etc. [14] with the Haar wavelet and Daubechies wavelet being used in this work. After transforming a time series, coefficients which ‘contain less information’ may be eliminated (*shrinkage*). This is achieved here by using the variance of the coefficients as a measure of information [14]. When combined with a neural network the overall model is known as a WP-MLP as shown below in Figure 2.

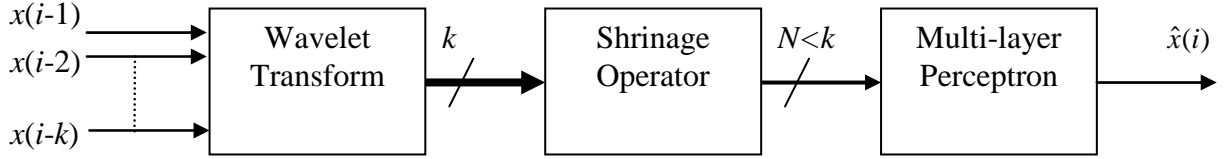


Figure 2 : WP-MLP Architecture

2.2 Neural Network Based Statistical Modelling

Several researchers have developed complex models and performed empirical studies of network jitter including those in [15] [16]. These studies show that network jitter follows a Laplacian distribution or a Normal distribution. The probability density function of the normal distribution is

$$f(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (4)$$

where σ is the standard deviation, and, μ is the mean.

The cumulative distribution function of the normal distribution is

$$F(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{(u-\mu)^2}{2\sigma^2}\right) du = \frac{1}{2} \left[1 + \operatorname{erf}\left(\frac{x-\mu}{\sigma \sqrt{2}}\right)\right] \quad (5)$$

In the proposed technique a neural network is used to predict the mean and variance parameters of a normal distribution model which is then used to calculate the desired playout delay value (ted), as shown in figure 3.

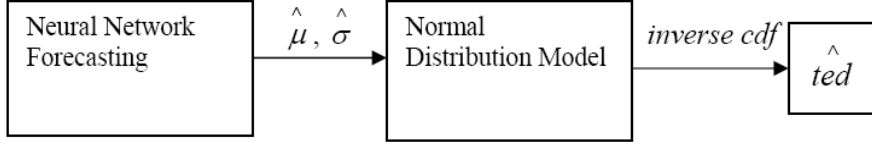


Figure 3 : Statistical-MLP Modelling

For a chosen mlp (maximum late packet loss percentage desired by user/application) value, there is a corresponding ted delay with any packets experiencing a jitter greater than ted being late for playout and hence discarded. For playout delay adaptation, the ted is chosen as the value of buffering delay that satisfies the condition $1 - cdf(ted) = mlp$.

$$\frac{1}{2} [1 + \operatorname{erf}(\frac{ted}{\sqrt{2}})] = 1 - mlp \quad (6)$$

$$ted = \sqrt{2} \times F^{-1}[2 \times (1 - mlp) - 1] \quad (7)$$

where $F^{-1}(x)$ refers to the inverse erf function.

3. Evaluation Methodology

In this paper, the various models were evaluated using real VoIP traces which were gathered using PJSIP [17], an open source VoIP application written in C, was adapted to measure the network jitter between two hosts. The application used in this paper first encodes the audio stream using G.729 B [18] into 20ms packets of length 80 bytes. Real Time Transport Protocol (RTP) is then used to sequence the packets and these are then encapsulated into a UDP packet for transmission across the internet. Since it was not feasible to take traces using terminals whose clocks were accurately synchronised, only information concerning inter-packet arrival times was available for these traces.

Several traces on international VoIP connections were taken ranging in duration from 5 to 10 hours of continuous duplex transmission from NUI, Galway to Tokyo (trace 1, a sample of which shown in Figure 4 below), Sydney (trace 2) and Chengdu (trace 3).

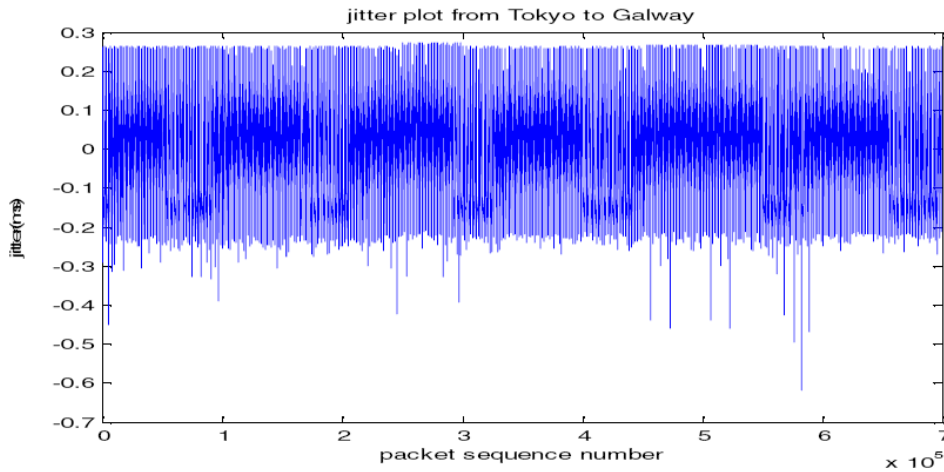


Figure 4 : Sample Jitter Plot (NUI, Galway to Tokyo)

At the receiver, an estimate of buffering delay (defined below) is used to allocate the playout time for each talkspurt as:

$$p^k(1) = r^k(1) + \hat{ted} \quad (8)$$

$$p^k(i) = p^k(1) + (i-1) \times 0.02 \quad i \neq 1 \quad (9)$$

where $p^k(1)$ is the playout time for first packet of the k^{th} talkspurt, $r^k(1)$ is the arriving time for first packet of the k^{th} talkspurt, \hat{ted} is the *estimated buffering delay* of the k^{th} talkspurt, $p^k(i)$ is the playout time for packet i of the k^{th} talkspurt and 0.02 seconds is ideal interval of the packet playout. \hat{ted} is the *estimated* according to the relative arriving time jitter ∇j .

The relative arriving time jitter of packet i $\nabla j(i)$ is defined as:

$$\nabla j^k(i) = r^k(i) - T^k(i) \quad (10)$$

$$\begin{aligned} T^k(i) &= r^k(1) + (i-1) \times 0. \\ (i \neq 1 \text{ Assume } T^k(1) &= r^k(1)) \end{aligned} \quad (11)$$

where $\nabla j^k(i)$ is the relative jitter for packet i of the k^{th} talkspurt, $r^k(i)$ is the arrival time packet i of the k^{th} talkspurt, $T^k(i)$ is the ideal arrival time packet i of the k^{th} talkspurt with no jitter.

Packets that arrive before their playout time slot ($p^k(i)$) are decoded using G.729 B. Packets that fail to arrive on time or that are dropped are ignored and are decoded instead using the G729 embedded Packet Loss Concealment (PLC) algorithm [19]. This algorithm attempts to interpolate the speech signal using previous packets in the stream.

The performance of each proposed model has been analyzed by three metrics:

1. Packet loss rate (the ratio of packets received, prior to $p^k(i)$, to those sent),
2. PESQ MOS metric, and
3. Additional buffering delay:

$$pd(i) = p^k(i) - r^k(i) \quad (12)$$

Perceptual evaluation of speech quality (PESQ) is a standard to measure the voice quality as published by the ITU-T. It compares a degraded speech signal, which is reconstructed after the network transmission and decoding, to an original signal and a MOS (mean opinion score) value is then produced. Commonly, the MOS value ranges from 0.0 (worst) to 4.5 (best) [20]. The overall algorithm evaluation that was used in the research scheme is shown below in Figure 5.

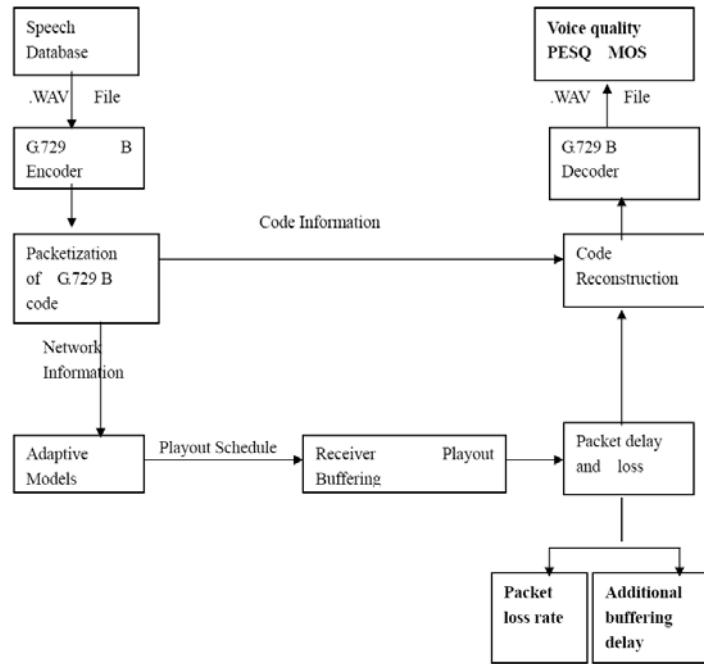


Figure 5 : Block Diagram for Performance Analysis Methodology

4. Results

When critically evaluating the performance of a playout scheduling algorithm, it is essential to consider both the additional delay and packet loss rate as shown in figure 6. This figure illustrates clearly the trade-off between additional delay and packet loss rates for the four different methods being proposed. The WP-MLP Haar based algorithm performs best in terms of packet loss (less than the limit for interactive audio, 5% [3]) and additional playout delay up to 400ms [2], which has been improved compared with the traditional MLP. The MLP also shows a good performance, compared with other methods. Comparatively, the results indicate that the RMLP based approach is not very suitable to be used in adaptive playout delay estimation. Alternatively, the Statistical-MLP model also shows a good performance and is very close to the WP-MLP Haar in terms of its abilities.

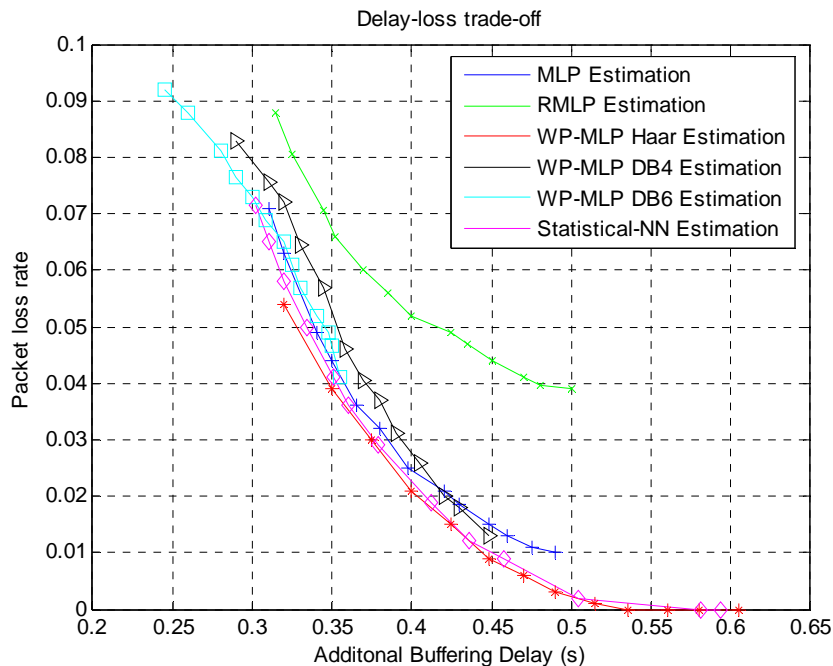


Figure 6 : Trade-off between Additional Buffering Delay and Packet Loss Rate

4.1 Results of PESQ MOS based Analysis

Table 2 below, gives a summary of the PESQ MOS score for the five techniques which were evaluated. These PESQ MOS scores were calculated by different lost packet information with the same additional buffering delay of 0.35s. The results show that the WP-MLP using the Haar wavelet also achieved the best performance when considering this perceptual based metric.

Model	WP-MLP Haar	Statistical-NN	MLP	WP-MLP DB4	RMLP	WP-MLP DB6
PESQ MOS	2.41234	2.40981	2.38173	2.09242	1.40953	2.38255

Table 2 : Comparison of Algorithm Performance using PESQ MOS Metric

5. Conclusions and Future Work

In this paper, several adaptive playout algorithms based on neural network have been presented and their performances evaluated. The effectiveness of these algorithms is evaluated using recorded traces by comparing the buffering delay and the packet loss ratios of each technique. Simulation results indicate that a Haar-Wavelets-Packet MLP adaptive scheduling scheme offers the best performance and flexibility for the process of adaptive playout delay estimation. A Statistical Modelling-MLP also shows a good performance very close to that of the WP-MLP (Haar). The WP-MLP DB4 and DB8 models also show an ability to minimise additional buffering delay but at the expense of higher packet loss rates. Future work will focus on the potential for improving the performance of the prediction performance of the WP-MLP by use of different mother wavelets and different levels of decomposition and the statistical-modelling Neural Network which may be improved by considering different combinations of statistical models and neural networks.

References

- [1] Davidson, J., Peters, J. (2000). Voice over IP Fundamentals. *Cisco Press*.
- [2] Telecommunication Standardization Sector of ITU. (1993) ITU-T Recommendation G.114. Technical report. *International Telecommunication Union*.
- [3] Jayant, N. S. (1980). Effects of packet loss on waveform coded speech. *Fifth Int. Conference on Computer Communications, Atlanta, Ga.*, 275–280.
- [4] Alvarez-Cuevas, F., Bertran, M., Oller, F., and Selga, J. (1993). Voice Synchronization in Packet Switching Networks. *IEEE Network Magazine*, 7:20–25.
- [5] Ramjee, R., Kurose, I., Towsley, D., and Schulzrinne, H. (1994). Adaptive playout mechanisms for packetized audio applications in wide-area networks. *IEEE INFOCOM*, 680–688.
- [6] Ramos, V., Barakat, C., and Altman, E. (2003). A moving average predictor for playout delay control in VoIP. *Quality of Service – IWQoS 2003, 11th International Workshop*, 155–173.
- [7] Liang, Y. J., Farber, N., and Girod, B. (2003). Adaptive playout scheduling and loss concealment for voice communications over IP networks. *IEEE Trans. Multimedia*, 5:532–543.
- [8] Moon, S. B., Kurose, I., and Towsley, D. (1998). Packet audio playout delay adjustment: Performance bounds and algorithms. *ACM Multimedia Systems*, 6:17-28.
- [9] Pinto, J., and Christensen, K. J. (1999). An algorithm for playout of packet voice based on adaptive adjustment of talkspurt silence periods. *IEEE Conf Local Computer Networks* (Lowell, MA), 224-231.
- [10] Agrawal, P., Chen, I. C., and Sreenan, C. J. (1998). Use of statistical methods to reduce delays for media playback buffering. *IEEE Int. Conf Multimedia Computing and Systems* (Austin, TX), 259-263.
- [11] Haykin, S. (1999). *Neural Networks: A Comprehensive Foundation*. Prentice. Hall, Upper Saddle River, NJ.
- [12] Zhang Q., Benveniste, A. (1992). Wavelet networks. *IEEE Trans. Neural Networks*, 3:889–898.
- [13] Hagan, M. T., Menhaj, M. B. (1994). Training Feedforward Networks with the Marquardt Algorithm. *IEEE Transactions on Neural Networks*, 5:989 -993.

- [14] Percival, D. B., Walden, A. T., (2000). *Wavelet Methods for Time Series Analysis*, Cambridge Univ. Press, Cambridge.
- [15] Zheng, L., Zhang, L., Xu, D. (2001). Characteristics of Network Delay and Delay Jitter and its Effect on Voice over IP (VoIP). *IEEE International Conference on Communications ICC 2001*, 1:122-126.
- [16] Li, M. P., Wilstrup, J., Jessen, R. and Petrich, D., (1999). A new method for jitter decomposition through its distribution tail fitting. *ITC Proceeding*, 788-794.
- [17] PJSIP Homepage, <http://www.pjsip.org/>
- [18] ITU-T Recommendation G.729 Annex B. (1996). A silence compression scheme for G.729 optimized for terminals conforming to Recommendation V.70.
- [19] ANSI. (2000). Packet Loss Concealment for use with ITU-T Recommendation G.711. *ANSI Recommendation T1.521a-2000 (Annex B)*.
- [20] Perceptual evaluation of speech quality (PESQ), an objective method for end-to-end speech quality assessment of narrowband telephone networks and speech codecs. (2001). *International Telecommunication Union*.