# PREDICTING QOS PARAMETERS FOR ATM TRAFFIC USING SHAPE-FUNCTION ESTIMATION.

Cormac Walsh[1] and N.G. Duffield[2]

## 1 Introduction and Summary

This paper addresses the provisioning of service capacity and buffer space in an ATM multiplexer of a large number of VBR sources. The aim is to guarantee sufficient quality of service (QoS), specified here in terms of the cell loss ratio (CLR), while at the same time making maximal use of system resources. Multiplexing gain is available in shared resources due to the statistical properties of the individual traffic streams which share the resources. For example, if one doubles the number of (identical) sources to be multiplexed, one need not, generally, double the rate and buffer size in order to maintain the same CLR.

Our method is to use empirical traffic statistics, which could be measured online, in order to predict these multiplexing gains. The prediction is done on the basis of the following scaling behaviour of queue-tail probabilities which has been shown to hold for a very general class of traffic models [1, 2, 4]. For an infinite buffer fed by $N$ independent identical sources served at rate $Ns$, the frequency with which the queue $Q^N$ exceeds a level $Nb$ satisfies the logarithmic asymptotic

$$\log P[Q^N > Nb] \sim -NI(b), \qquad \text{as } N \to \infty, \tag{1}$$

where the *shape function* $I(b)$ depends on $s$ and the detailed traffic characteristics. Although this estimate is formulated for the asymptotic behaviour of tail probabilities in an infinite buffer, it is conservative: it gives an upper bound on the CLR from a buffer of size $Nb$. (Actually, in some cases the difference itself can be estimated). Moreover, as we shall discuss, corrections for finite $N$ are available (see [3, 9]). Another useful property is that (1) does not require that $b$, the buffer allocation per source, be large. Thus it can be used to describe cell-level, as well as burst-level, queueing behaviour. In this, it is distinguished from the large body of results about asymptotic behaviour of tail probabilities for large $b$, and the consequent effective bandwidth approximation (see [6, 10] and references therein).

The format of the paper is as follows: in section 2, we describe the basis for (1); in section 3, we describe an estimator of the shape function $I$; in section 4, we review a scheme proposed by Courcoubetis, Fouskas and Weber [3] to estimate the CLR which makes use of finite $N$ corrections to scaling behaviour of the kind shown in (1); in sections 5 and 6, we compare the accuracy and the reliability of the two by applying them to both simulated and to real ATM traffic.

## 2 The Shape Function

We model an ATM switch as a single server queue with $N$ (possibly heterogeneous) stationary and ergodic arrival streams, $\{A_t^{(1)}, \ldots, A_t^{(N)}\}$. When we scale $N$, we shall assume the sources are served at constant load so that the service rate is $sN$ for some constant $s$. Define the finite-time *cumulant generating function* (finite-time CGF) for each source $i$ at each time $t$ as

$$\lambda_t^{(i)}(\theta) := \frac{1}{t} \log E e^{\theta A_t^{(i)}} \tag{2}$$

The finite-time CGFs are related to the *effective bandwidth* of the source: $\alpha^{(i)}(\theta) = \theta^{-1} \lim_{t \to \infty} \lambda_t^{(i)}(\theta)$.

[1]School of Mathematical Sciences, Dublin City University, Dublin 9, Ireland and Dublin Institute for Advanced Studies, 10 Burlington Road, Dublin 4, Ireland; E-mail `walsh@stp.dias.ie`

[2]Room 2C-323; AT&T Laboratories, 600 Mountain Avenue, Murray Hill, NJ 07974, USA; E-mail `duffield@research.att.com`

We make the assumption that the finite-time CGF for the combined traffic stream exists in the limit as $N \to \infty$:

$$\lambda_t(\theta) := \lim_{N \to \infty} \frac{1}{Nt} \log E e^{\theta \sum_i A_t^{(i)}}.$$

Furthermore, we require that this limit exists uniformly for all $t$ sufficiently large. This assumption is satisfied by i.i.d. superpositions and also by heterogeneous superpositions where the proportion of each type of source is held constant. In the case of independent heterogeneous superpositions of $J$ types of sources indexed by $j \in \{1, \ldots, J\}$, we have that $\lambda_t(\theta) = \sum p_j \lambda_t^{(j)}(\theta)$ where $p_j$ is the proportion of sources of type $j$. Under these assumptions, (1) holds with the shape function $I$ given by

$$I(b) = \inf_{t \geq 0} (t\lambda_t)^*(b + st) \tag{3}$$

where $f^*(x) := \sup_{y \in \mathbf{R}} \{xy - f(y)\}$ is the *Legendre-Fenchel transform* of the function $f$; see [1, 2, 4]. The time $\tau = \arg\inf_{t \geq 0} (t\lambda_t)^*(b)$ at which the infimum above is attained may be interpreted as the most likely timescale on which the buffer overflows.

## 3 The Shape Function Estimator (SFE).

Given a sample realization $\{X_1, X_2, \ldots\}$ of a traffic stream we may estimate the finite-time CGFs of its source as follows. First form all blocks of length $t$:

$$\tilde{X}_1 := \sum_{i=1}^{t} X_i, \qquad \tilde{X}_2 := \sum_{i=2}^{t+1} X_i, \qquad \ldots$$

Assuming stationarity of the arrival process, we use these overlapping blocks to get an estimate of $\lambda_t$ by replacing the expectation in (2) with an empirical mean:

$$\hat{\lambda}_t(\theta) := \frac{1}{t} \log \frac{1}{K} \sum_{k=1}^{K} e^{\theta \tilde{X}_k}$$

where $K$ is the number of blocks formed. We assume that the sources are independent and so we may combine the estimates to form an estimate of $\lambda_t$:

$$\hat{\lambda}_t(\theta) := \frac{1}{N} \sum_{j=1}^{N} \hat{\lambda}_t^{(j)}$$

Then for each $t$ we merely perform the infimum and Legendre transform, mirroring (3), in order to form the estimate

$$\hat{I}(b) := \inf_{t \geq 0} (t\hat{\lambda}_t)^*(b + st)$$

It is worth distinguishing this procedure from that of [5] where, for some (large) $T$, $\hat{\lambda}_T$ was used to estimate the limiting CGF, $\lambda(\theta) = \lim_{t \to \infty} \lambda_t(\theta)$, and hence estimate the (large buffer asymptotic) effective bandwidth by $\hat{\alpha}(\theta) := \theta^{-1} \hat{\lambda}_T(\theta)$. There is a tradeoff here in choosing a value for $T$: too small a value and $\lambda_T(\theta)$ will not be close enough to $\lambda(\theta)$, too large a value and the variance of the estimator will be large (see [7]). Also, it is difficult to automate the choice of block size. However, we shall see that the shape-function estimator does not suffer from this problem. Once we have an estimate of the shape-function, we may use it to predict cell loss ratios; the CLR is estimated by the *shape function estimator* (SFE):

$$CLR(b) \approx e^{-N\hat{I}(b)}$$

# 4 The Empirical Loss Rate Estimator (ELRE)

Courcoubetis *et. al.* [3] propose an estimator of the CLR which is also based on the large $N$ asymptotics of the system. The ELRE essentially involves the direct estimation of the cell loss rate for subsets of the $N$ traffic sources by simulating the queue with the appropriately scaled buffer size and service rate. Denote the number of sources in the $i^{th}$ subset by $N_i$ and the observed CLR of this subset by $\Phi_i$. Then parameters $A$,$B$, and $C$ are chosen so as to best fit a curve of the form $A + B \log N_i + CN_i$ to the graph of $\log \Phi_i$ vs. $N_i$. The CLR, $\Phi(N)$, may then be calculated for any $N$:

$$\log \Phi(N) = A + B \log N + CN$$

The parameter $C$ corresponds to $-\hat{I}(b)$ above, so we may compare these estimates directly. The $B \log N$ term is motivated by the Bahadur-Rao refinement to large deviation asymptotics for sums of independent random variables; see also [9]. The authors of [3] recommend using non-overlapping subsets of sources (each source is included in no more than one subset). They do this because they want the CLR estimates to be independent in order to facilitate a $\chi^2$ goodness-of-fit test.

# 5 Comparisons for Model-based Simulations

We compare the estimators' prediction of the CLR and of the shape function for model traffic. The traffic model we use is the two state discrete time Markov model. This model produces a cell every clock cycle while in state 1 and no cells while in state 0. The transition probabilities are $P(X_{n+1} = 1|X_n = 0) = 1 - P(X_{n+1} = 0|X_n = 0) = 1/16$ and $P(X_{n+1} = 0|X_n = 1) = 1 - P(X_{n+1} = 1|X_n = 1) = 3/16$. With these parameters the model is positively correlated; cells tend to arrive in bursts, in this case of mean length 16/3 clock cycles.

We generated 100 different sets of samples, each of 20 independent sources. The left hand side of figure 1 shows the interquartile range of the predictions of the CLR for 30 sources based on measurements of 20 sources. The ELRE does better than the SFE; as we might expect from the discussion in the introduction, the SFE is conservative in estimating the CLR. However, we can also calculate the shape function $I$ numerically for the model and compare it with the estimates of it obtained using both schemes ($\hat{I}$ in the SFE, $-C$ in the ELRE). The quartiles are shown on the right hand side of figure 1. In each plot, the horizontal line represents the true value of the shape function evaluated at the relevant buffer size and service rate. The estimates from the SFE are seen to be tightly clustered about the true value, their spread decreasing as the sample size is increased. The estimates from the ELRE also improve with increasing sample size. However, their spread is always greater than those from the shape function estimator and, furthermore, they exhibit bias. This bias is not consistent: the shape function may be over or underestimated at different $b$. We found the bias to depend critically on the range of source numbers over which the $\log \Phi$ vs. $N$ curve was fitted; a different range was used for each sample size and so the bias changed accordingly.

Obtaining an accurate estimate of $I$ is particularly important in applications where extrapolations are made to values of $N$ which are far greater than the number of samples used; this is because, even for the ELRE, the shape function estimate $C$ becomes dominant for large $N$. These considerations may be relevant in provisioning capacity of a large number for VBR sources on the basis of observed characteristics of a small set of source models or traces. However, for moderate $N$, the corrections $A + B \log N$ in the ELRE give better predictions. For this reason we are now experimenting with a combined estimator that fits a two parameter curve of the form $A + B \log N - \hat{I}(b)N$ to the graph of $\log \Phi$ vs. $N$ where $\hat{I}(b)$ is now a constant determined by the shape function estimator. This estimator should provide a narrower spread of estimates than the empirical loss rate estimator while avoiding the bias of the pure shape function estimator.

# 6 Trace driven simulations

To conclude, we shall apply the estimators to real ATM traffic. We chose to use the well known 'Starwars' video data set produced by Garrett and Vetterli [8]. Using the number of cells per slice

(there are 30 slices per frame so the slice time is approximately 1.4ms), we split the traffic into 100 disjoint segments of equal length and multiplexed them together. We then applied the estimators to the resulting traffic stream as we did for the traffic produced from the model. Note that we are not testing predictions of future cell loss here, we are merely checking to see if the two estimators give us a consistent picture of the shape function of this traffic set. Figure 2 shows both estimates of the shape function plotted against buffer size for a particular service rate. It is seen that the two estimators agree reasonably well especially at moderately large buffer sizes.

# References

[1] D.D. Botvich and N.G. Duffield (1995). Large deviations, economies of scale, and the shape of the loss curve in large multiplexers. *Queueing Systems.* **20:** 293-320.

[2] C. Courcoubetis and R. Weber (1996). Buffer overflow asymptotics for a switch handling many traffic sources. *J. Appl. Prob.* **33:** 886–903

[3] C. Courcoubetis, G. Fouskas and R. Weber (1995). An on-line estimation procedure for cell-loss probabilities in ATM links. 3rd IFIP Workshop on Performance, Modeling and Evaluation of ATM Networks. UK, July 1995.

[4] N.G. Duffield (1996). Economies of scale in queues with sources having power-law large deviation scalings. *J. Appl. Prob.,* **33:** 840–857.

[5] N.G. Duffield, J.T. Lewis, N. O'Connell, R. Russell and F. Toomey (1995). Entropy of ATM traffic streams: a tool for estimating QoS parameters. *IEEE J. Selected Areas in Commun.* **13** 981–990.

[6] N.G. Duffield and N. O'Connell (1995). Large deviations and overflow probabilities for the general single-server queue, with applications, *Math. Proc. Cam. Phil. Soc.,* **118:**363–374.

[7] A. J. Ganesh (1996). Bias Correction in Effective Bandwidth Estimation, Computer Systems Group Report, ECS-CSG-23-96, May 1996.

[8] M.W. Garrett and M. Vetterli (1991). Congestion control strategies for packet video. In: *Proceedings of Fourth International Workshop on Packet Video*, Kyoto, Japan, August 1991.

[9] M. Montgomery and G. de Veciana (1996) On the relevance of time-scales in performance oriented traffic characterizations. Proceedings IEEE INFOCOM'96, pp513–520.

[10] W. Whitt (1993). Tail probabilities with statistical multiplexing and effective bandwidths in multi-class queues, Telecom. Syst. **2:** 71-107
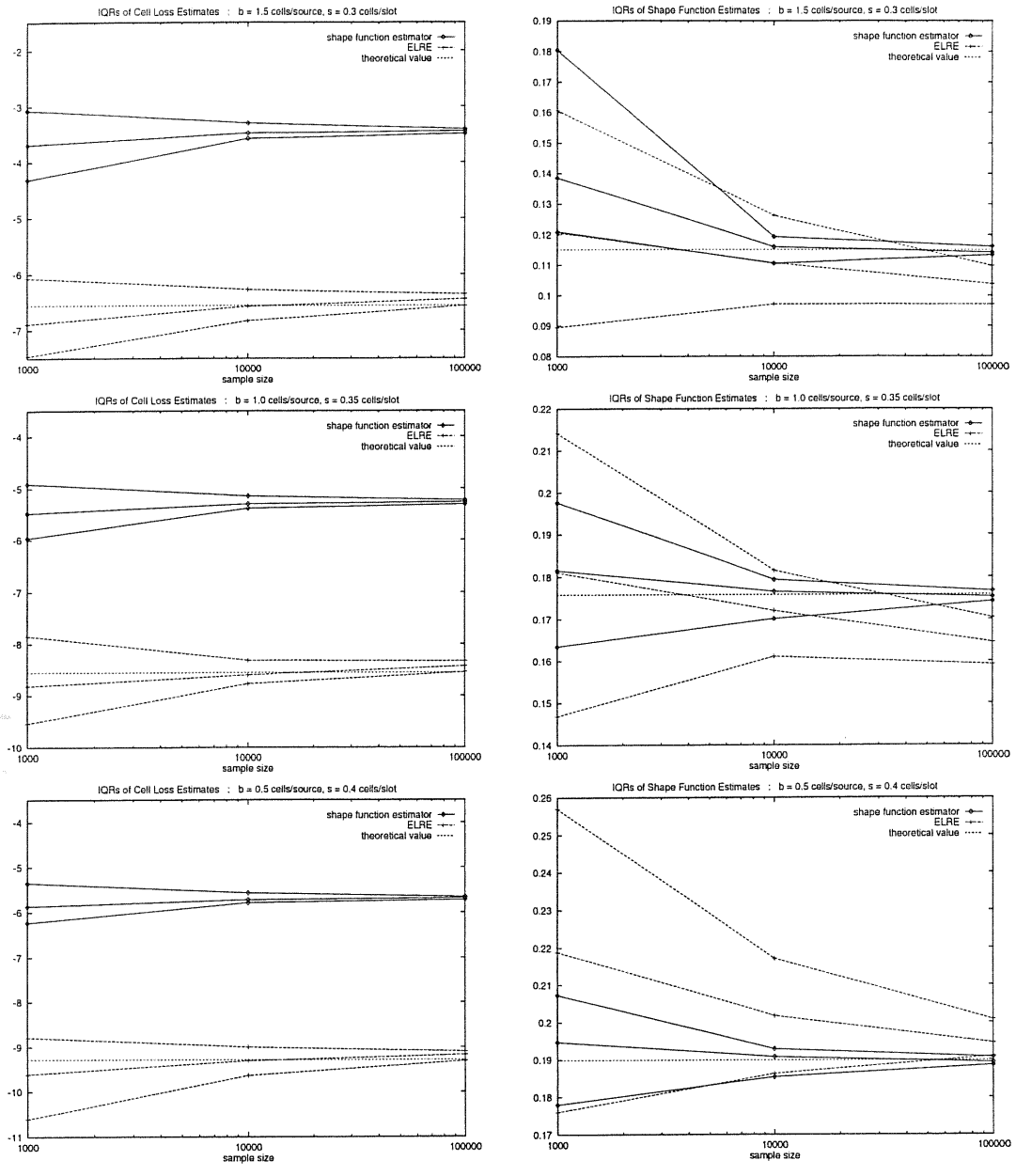
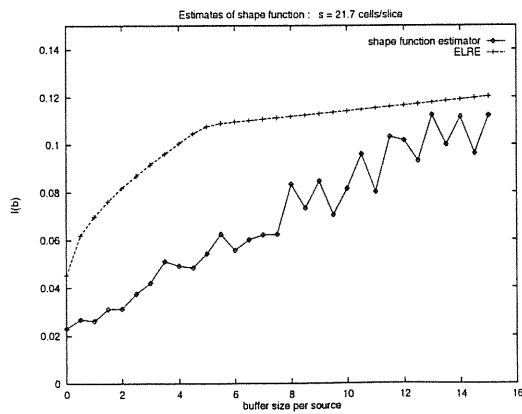Figure 1: Interquartile range of loss estimates and shape function estimates



Figure 2: Estimates of the shape function of "StarWars"