

# Sensitivity Analysis for Discrete-Time Randomized Service Priority Queues\*

George Kesidis<sup>1</sup> and Takis Konstantopoulos<sup>2</sup> and Michael Zazanis<sup>3</sup>  
Proc. IEEE CDC'94, Orlando, FL, pp. 2627-2630, Dec. 1994.

February 23, 1995

## Abstract

We consider a collection of queues operating in parallel and sharing a common server via the idling randomized service priority sharing rule. We assume that the service times of the customers (cells) are all constant and identical. This paper is concerned with estimating the sensitivity of the tail of the distribution of the typical sojourn time through these queues to the fraction of server bandwidth given to them by randomized service sharing. Two approaches are considered: "smoothed" perturbation analysis and stochastic intensity-based estimators.

## 1 Introduction

Asynchronous Transfer Mode (ATM) is the emerging standard for high-speed, integrated networks [2]. ATM networks will handle different kinds of traffic, consisting of streams of small fixed-length packets called cells, having different performance requirements from the network. For instance, a voice connection will require a small end-to-end delay and can tolerate significant packet loss. A data connection, however, may require an extremely small packet loss probability and may tolerate significant delay [7]. In this paper, we focus on the buffer design wherein each traffic class occupies a separate buffer and the buffers share a server via a "randomized" service priority sharing rule [8].

We now describe the operation of the idling, non-interrupting, randomized service sharing policy on a group of  $K$  queues with independent, stationary and ergodic sources. The server is assigned to a queue until it finishes serving a single cell (non-interrupting). If

the server is assigned to an empty queue, it waits there for an amount of time equal to the service time of one cell (idling). The assignment process is determined by i.i.d. random variables  $\{\gamma_n\}$  that are uniformly distributed over  $[0, 1]$ . The interval  $[0, 1]$  is partitioned into  $K$  smaller intervals  $\{A_i\}_{i=1}^K$  so that if  $\gamma_n \in A_i$ , then the server will be assigned to the  $i^{\text{th}}$  queue in the  $n^{\text{th}}$  service epoch.

Let  $\theta = P(\gamma_n \in A_i)$  be the fraction of service bandwidth the  $i^{\text{th}}$  queue receives on average. In this paper, we are concerned with estimating the sensitivity

$$\eta := \frac{d}{d\theta} P(T(\theta) \geq B) \quad (1)$$

where  $T(\theta)$  is the sojourn time of a typical cell in this queue (to ease the notation we omit the index  $i$ ). The choice of  $\theta$  clearly depends on performance requirements of the traffic using this queue. For reliable network design, the choice of  $\theta$  should depend on  $\eta$  as well.

The paper is organized as follows. In Section 2, we describe a smoothed perturbation analysis estimator for  $\eta$ . A recursive update formula to estimate  $\eta$  from a single sample path is presented in Section 3. In Section 4, we present an alternative SPA formula which can be directly implemented without the need for a recursive procedure. Finally, an approximate method based on stochastic intensities is described in Section 5 and conclusions are drawn in Section 6.

## 2 Smoothed Perturbation Analysis

Consider a "G/D/RS" queue; i.e., a queue with a discrete-time stationary and ergodic cell arrival process that shares the server with other such queues according to the randomized service priority discipline. Let  $\theta$  be the fraction of service bandwidth

---

\*Supported by NSERC of Canada and NSF Research Initiation Award NCR-9211343, 1. Elec. & Comp. Eng. Dept., University of Waterloo, Waterloo, ON, Canada, N2L 3G1, 2. Dept. of Elec. & Comp. Eng., University of Texas, Austin TX 78712, 3. Dept. of Industrial Eng., University of Massachusetts, Amherst MA 01003.

this queue receives on average. We want to find  $dP(T(\theta) \geq B)/d\theta$  where  $T(\theta)$  is the sojourn time of a typical cell.

Assume that this queue is simulated for  $N$  cell departures and to the  $n^{\text{th}}$  arrival ( $n = 1, \dots, N$ ) associate an infinite sequence of i.i.d. uniform  $[0, 1]$  random variables  $\{\xi_i^n\}_{i=1}^\infty$ . The amount of “virtual service” required by the  $n^{\text{th}}$  cell is

$$\sigma_n(\theta) = \inf\{i : \xi_i^n \leq \theta\}. \quad (2)$$

Note that  $\sigma_n(\theta)$  has a geometric distribution. Define the random vector  $\sigma(\theta) := (\sigma_1(\theta), \dots, \sigma_N(\theta))$ . Now consider two simulations of the G/D/RS queue: one using the parameter  $\theta$  and the other using  $\theta - \Delta\theta$  where  $0 < \Delta\theta \ll 1$ . Only the  $\theta$ -simulation will actually be conducted. Below we describe how to estimate  $dP(T(\theta) \geq B)/d\theta$  given the results of the  $\theta$ -simulation alone.

Using the notation of Suri [9], for the  $\theta$ -simulation of length  $N$  cell departures let:  $BP_j$  be the  $j^{\text{th}}$  busy period,  $t_j$  be the starting time of  $BP_j$  = the  $j^{\text{th}}$  arrival time of a cell to an empty queue,  $s_j$  be the ending time of  $BP_j$  = the  $j^{\text{th}}$  departure time of a cell leaving the queue empty,  $M$  be the number of  $BP$ 's for the  $N$  arrivals simulated,  $k_{j+1}$  be the index of the last cell of  $BP_j$  with  $k_1 = 0$ ,  $\delta_j = t_j - s_{j-1}$ ,  $\Delta_j^i = \sum_{k=j+1}^i \delta_k$  when  $i > j$  and  $\Delta_j^i = 0$  when  $i = j$ .

Define  $\mathbf{A} = \{x \in \mathbf{R}^N \mid x = k\mathbf{e}_j, k = 0, 1, 2, \dots, j = 1, 2, 3, \dots\}$  where  $\mathbf{e}_j$  is the unit vector whose  $j^{\text{th}}$  entry is 1 and all other entries are 0. Note that

$$\begin{aligned} P(\sigma(\theta - \Delta\theta) - \sigma(\theta) = 0 \mid \sigma(\theta)) &= (1 - \Delta\theta/\theta)^N = 1 + o(\Delta\theta), \\ P(\sigma(\theta - \Delta\theta) - \sigma(\theta) = k\mathbf{e}_j \mid \sigma(\theta)) &= (\theta - \Delta\theta)(1 - \theta + \Delta\theta)^{k-1} \frac{\Delta\theta}{\theta} (1 - \frac{\Delta\theta}{\theta})^{N-1} \\ &= (1 - \theta)^{k-1} \Delta\theta + o(\Delta\theta), \\ P(\sigma(\theta - \Delta\theta) - \sigma(\theta) \notin \mathbf{A} \mid \sigma(\theta)) &= o(\Delta\theta). \end{aligned}$$

Thus, to find an expression for an estimate of

$$\begin{aligned} \frac{d}{d\theta} P(T(\theta) \geq B) &= \lim_{\Delta\theta \rightarrow 0} \frac{1}{\Delta\theta} P(T(\theta) \geq B) - P(T(\theta - \Delta\theta) \geq B) \\ &= \lim_{\Delta\theta \rightarrow 0} \frac{1}{\Delta\theta} E(E[\mathbf{1}\{T(\theta) \geq B\} - \mathbf{1}\{T(\theta - \Delta\theta) \geq B\} \mid \sigma(\theta), \text{Arrivals}]), \end{aligned}$$

we need only consider the case where a single cell's virtual service extends when  $\theta \rightarrow \theta - \Delta\theta$ . Note that  $E(\mathbf{1}\{T(\theta) \geq B\} \mid \sigma(\theta), \text{Arrivals}) = E(\mathbf{1}\{T(\theta) \geq B\}) = P(T(\theta) \geq B)$ .

By an argument similar to that in [9], Section VI.B (the “smoothed perturbation analysis” of Gong and Ho [4]), an estimate of  $dP(T(\theta) \geq B)/d\theta$  is

$$\hat{\eta}(N) := -\frac{1}{N} \sum_{j=1}^M \sum_{n=k_j+1}^{k_{j+1}} \sum_{a=1}^{\infty} (1 - \theta)^{a-1} R(j, n, a) \quad (3)$$

where

$$\begin{aligned} R(j, n, a) &:= \sum_{b=j}^M \sum_{m=\max\{n, k_b+1\}}^{k_{b+1}} \mathbf{1}\{B - [a - \Delta_j^b]^+ \leq T_m < B\} \end{aligned} \quad (4)$$

where  $T_m$  is the sojourn time of the  $m^{\text{th}}$  cell in the  $\theta$ -simulation.

We now explain equations (3) and (4). Cell  $n$  residing in  $BP_j$  has its virtual service time extended by  $a$  units when  $\theta \rightarrow \theta - \Delta\theta$ . Consequently, some busy periods, beginning with  $BP_j$ , may coalesce. The result is that the sojourn times of some cells will increase when  $\theta \rightarrow \theta - \Delta\theta$ :

$$\begin{aligned} T_m(\theta - \Delta\theta) - T_m(\theta) &= \begin{cases} a, & n \leq m \leq k_{j+1} \\ (a - \Delta_j^{j+1})^+, & k_{j+1} + 1 \leq m \leq k_{j+2} \\ (a - \Delta_j^{j+2})^+, & k_{j+2} + 1 \leq m \leq k_{j+3} \\ \text{etc.} \end{cases} \end{aligned}$$

Equation (4) then follows from

$$\begin{aligned} \mathbf{1}\{T_m(\theta) \geq B\} - \mathbf{1}\{T_m(\theta - \Delta\theta) \geq B\} &= \mathbf{1}\{T_m(\theta) \geq B\} - \mathbf{1}\{T_m(\theta) + (a - \Delta_j^b)^+ \geq B\} \\ &= -\mathbf{1}\{B - (a - \Delta_j^b)^+ \leq T_m(\theta) < B\} \end{aligned}$$

### 3 Recursive Update Formula for the Estimate

In this section we find a simple formula for  $\delta(N + 1) := \hat{\eta}(N + 1) - \hat{\eta}(N)$ . This formula can be used to update our estimate of the sensitivity after every cell departure of the  $\theta$ -simulation.

In the following we assume that the random quantities  $k_j$  and  $M$  are evaluated for first  $N$  departed cells. Consider two cases for the cell  $N + 1$ :

If cell  $N + 1$  is a member of  $BP_M$ ,

$$\begin{aligned} \delta(N + 1) &= -\sum_{a=1}^{\infty} (1 - \theta)^{a-1} \mathbf{1}\{B - a \leq T_{N+1} < B\} \end{aligned}$$

$$\begin{aligned}
& - \sum_{j=1}^M \sum_{n=k_{j+1}}^{k_{j+1}} \sum_{a=1}^{\infty} (1-\theta)^{a-1} \times \\
& \quad \mathbf{1}\{B - (a - \Delta_j^M)^+ \leq T_{N+1} < B\} \\
& = - \frac{\mathbf{1}\{T_{N+1} < B\}}{\theta} [(1-\theta)^{B-T_{N+1}-1} + \\
& \quad \sum_{j=1}^M (k_{j+1} - k_j)(1-\theta)^{B-T_{N+1}+\Delta_j^M-1}].
\end{aligned}$$

Note that the first term in the equations above is due to the virtual service time of cell  $N+1$  increasing by  $a$ . The second term is due to the virtual service time of one of the first  $N$  cells increasing by  $a$  causing the queuing time of cell  $N+1$  to increase.

If cell  $N+1$  is not a member of  $BP_M$  (i.e., cell  $N+1$  begins the new busy period  $BP_{M+1}$ ),

$$\begin{aligned}
& \delta(N+1) \\
& = - \sum_{j=1}^{M+1} \sum_{n=k_{j+1}}^{k_{j+1}} \sum_{a=1}^{\infty} (1-\theta)^{a-1} \times \\
& \quad \mathbf{1}\{B - (a - \Delta_j^{M+1})^+ \leq T_{N+1} < B\} \\
& = - \frac{\mathbf{1}\{T_{N+1} < B\}}{\theta} \sum_{j=1}^{M+1} (k_{j+1} - k_j) \times \\
& \quad (1-\theta)^{B-T_{N+1}+\Delta_j^{M+1}-1}.
\end{aligned}$$

Note that  $k_{M+2} - k_{M+1} = 1$  and the two cases yield very similar expressions.

## 4 An Alternative SPA Estimator

An alternative estimator can be obtained if in the above SPA analysis the perturbation  $\Delta\theta$  is taken in the opposite direction. Denote by  $W_i(\theta)$  the workload process when the parameter value is equal to  $\theta$  (the *nominal* sample path). From (2) it follows that  $\sigma(\theta + \Delta\theta) \leq \sigma(\theta)$  w.p.1 and hence that the cells that arrive to an idle system in the nominal path  $i$  (“lucky” cells) will remain lucky in the perturbed path as well. To keep the notation consistent we will use  $P$  to designate the Palm probability w.r.t. the arrival process while  $P^*$  will denote the Palm probability w.r.t. the *lucky* arrivals at parameter value  $\theta$ . The cycle formula between these two measures gives

$$\begin{aligned}
& \frac{1}{\Delta\theta} [P(T_0(\theta + \Delta\theta) > B) - P(T_0(\theta) > B)] = \quad (5) \\
& \quad \frac{1}{E^*Q} E^* \left[ \sum_{i=0}^{Q-1} \mathbf{1}(T_i(\theta + \Delta\theta) > B) - \mathbf{1}(T_i(\theta) > B) \right],
\end{aligned}$$

where  $Q$  is the number of customers in the first busy period of the nominal sample path. In general terms the difference between this approach and that of Section 2 is that we now need only worry about a single busy period breaking up instead of several busy periods coalescing. To implement the SPA algorithm we will condition w.r.t  $\mathcal{F}$ , the whole history of the *nominal* sample path. An analysis similar to that of Section 2 gives

$$\begin{aligned}
& \frac{d}{d\theta} P(T_0(\theta) > B) = - \frac{1}{(1-\theta)E^*Q} \times \\
& \quad E^* \left[ \sum_{i=0}^{Q-1} \sum_{j=i}^{Q-1} \sum_{k=1}^{\sigma_{i-1}} \mathbf{1}(B < T_i(\theta) \leq B + L_{ij} \wedge k) \right],
\end{aligned}$$

where

$$L_{ij} = \min\{W_{i+1}, \dots, W_j\}, \quad j \geq i,$$

with the convention that the minimum element of the empty set is  $+\infty$ .

## 5 Stochastic Intensity Based Estimators

Consider again a G/D/RS queue but assume that the server has nonzero setup times. This means that it takes the server a nonzero amount of time to switch from one queue to another. The setup times are small compared to the service times (still taken to be of unit duration) but random. Let  $f_s(x)$  denote their common density, supported on the interval  $[0, \epsilon]$ , where  $\epsilon < 1$ . The virtual service time  $\sigma_n(\theta)$  of the  $n^{\text{th}}$  cell is now seen to have density

$$g(\theta, x) = \sum_{k=1}^{\infty} \theta(1-\theta)^{k-1} f^{(k)}(x), \quad (6)$$

where  $f(x) = f_s(1+x)$ , and  $f^{(k)}(x)$  is the  $k$ -fold convolution of  $f$  with itself. Let  $G(\theta, x) = \int_0^x g(\theta, y) dy$  be corresponding distribution function. The random variable  $\sigma_n(\theta)$  is now generated by the formula  $\sigma_n(\theta) = G^{-1}(\theta, \xi_n)$ , where  $G^{-1}$  is the inverse function of  $G$  with respect to the second variable, and  $\xi_n$  is a sequence of i.i.d. random variables, uniformly distributed in the interval  $[0, 1]$ . Finally, we let

$$\sigma'_n(\theta) = \frac{d}{d\theta} G^{-1}(\theta, \xi_n). \quad (7)$$

It can be seen that this derivative is defined (i.e., it is finite) for all  $\xi_n$  outside an interval of size of the

order of  $1 - (1 - \theta)^{1/\epsilon}$ . This is due to the nature of the density  $g$  defined in (6) as the sum of convolutions of a density  $f$  supported on an interval of size  $\epsilon$ . It is natural to require that this size be small so that the error in the algorithm to follow is negligible. We thus need  $(1 - \theta)^{1/\epsilon} \approx 1$ . This is for instance the case if  $\epsilon$  is large or if  $\theta$  is small. It is thus conjectured that the algorithm works well for low priority classes.

These assumptions, namely in cases where one can afford an additional randomization for the service times, lead to a considerable simplification of the perturbation analysis estimator. Indeed, it is shown in [6], that an infinitesimal perturbation analysis estimator can be constructed. The construction is based on knowledge of the stochastic intensity, say  $\alpha_t$ , of the arrival process of the queue under consideration. Recall that the stochastic intensity of a point process, c.f. [1], with respect to a  $\sigma$ -field  $\mathcal{F}_t$  of observations (here taken to be the information of the simulated sample path of the queue under consideration up to time  $t$ ), is defined by

$$\alpha_t = \lim_{\delta \rightarrow 0} \frac{1}{\delta} E[N(t, t + \delta) | \mathcal{F}_t],$$

where  $N(t, t + \delta)$  is the number of arrivals between  $t$  and  $t + \delta$ . Note that not every point process has a stochastic intensity, but most models encountered in practice do. For instance, when the arrival process is renewal, it has stochastic intensity (with respect to the  $\sigma$ -field of the sample path up to time  $t$ )  $\alpha_t = h(Z_t)$ . Here  $h$  is the hazard rate of the interarrival time ( $h(x)$  defined as  $f_a(x)/(1 - F_a(x))$ , with  $f_a$  being the density of the interarrival time and  $F_a$  its distribution function), and  $Z_t$  is the distance between time  $t$  and the previously observed arrival. Likewise, simple formulas for the stochastic intensity can be found for most models used in practice (e.g., Markov modulated Poisson processes).

Let  $W_t(\theta)$  be the total work in the queue at time  $t$ , as accounted for by the (remaining) virtual service times of the cells in the queue. Let now  $W(\theta)$  be the total queuing delay of a typical cell in steady state. This is related to the sojourn time  $T(\theta)$  by  $T(\theta) = W(\theta) + \sigma_0(\theta)$ , where  $\sigma_0(\theta)$  is a typical virtual service time with density  $g$  as above. We now simulate the queue for a total of  $N$  cell arrivals (say  $T_N$  is the time of the  $N^{\text{th}}$  arrival) and define the following quantities. Let  $D_t(\theta)$  be the sum of the derivatives of the service times, given by (7), of all cells from the start of the busy period that contains  $t$ , up to the last cell arriving before  $t$ . If  $t$  is in an idle period then  $D_t(\theta)$  is taken to be zero. We are interested in estimating the derivative of  $P(W(\theta) > x)$ . The queue is simulated for  $N$  cell

arrivals, and observations are made at times at which the total work  $W_t(\theta)$  *downcrosses* level  $x$ . Call  $S_j$  these times. It is shown in [6] that an estimator for  $\frac{d}{d\theta} P(W(\theta) \geq x)$  is given by the quantity

$$\zeta(N) := \frac{1}{N} \sum_{j \geq 1, S_j < T_N} D_{S_j}(\theta) \alpha_{S_j}.$$

In other words, at each downcrossing time  $S_j$ ,  $j = 1, 2, \dots$ , the stochastic intensity  $\alpha_{S_j}$  is computed and multiplied by the accumulated perturbation  $D_{S_j}(\theta)$ . The sum of these products up to the  $N^{\text{th}}$  arrival divided by  $N$  gives a simple expression for the sensitivity estimator. The reader is referred to [6] for the relevant details and proofs.

## 6 Conclusions

We have described perturbation analysis estimators for the sensitivity of tail of the queuing delay distribution for a queue sharing a server via idling randomized service. We assumed a stationary and ergodic source of customers (cells) all requiring the same deterministic amount of service. Currently, we are working on this problem for state-dependent (in particular, work conserving) randomized service disciplines.

## References

- [1] P. Brémaud. *Point Processes and Queues*. Springer, 1981.
- [2] J. Filipiak. *Real Time Network Management*. North-Holland, New York, NY, 1991.
- [3] W.B. Gong, Smoothed Perturbation Analysis Algorithm for a G/G/1 Routing Problem, in *Proc. of the 1988 Winter Simulation Conference*, M. Abrams, P. Haight, J. Comfort, Eds., 525-531, 1988.
- [4] W.B. Gong and Y.C. Ho, Smoothed (conditional) perturbation analysis of discrete event dynamic systems. *IEEE Trans. Auto. Control*, Vol. 32, No. 10: pp. 858-866, 1987.
- [5] T. Konstantopoulos and M. Zazanis. Sensitivity Analysis for Stationary and Ergodic Queues. *Adv. Appl. Prob.*, Vol. 24, 738-750, 1992.
- [6] T. Konstantopoulos and M. Zazanis. Stochastic Intensity Based Sensitivity Estimators for Stationary and Ergodic Queues. *Preprint*, Aug., 1993.

- [7] A. A. Lazar, G. Pacifici, and J. S. White. Real-time traffic measurement on MAGNET II. *IEEE JSAC*, Vol. 8, No. 3:467–493, April 1990.
- [8] J.M. Pitts and J.A. Schormans. Analysis of ATM switch model with time priorities. *Electronic Letters*, Vol. 26, No. 15:1192,1193, July 1990.
- [9] R. Suri. Perturbation analysis: the state of the art and research issues explained via the GI/G/1 queue. *Proceedings of the IEEE*, Vol. 77, no. 1:114–137, 1989.