

STOCHASTIC FLUID MODELS FOR CONTROL AND OPTIMIZATION OF SYSTEMS WITH QUALITY OF SERVICE REQUIREMENTS¹

Christos G. Cassandras, Gang Sun, and Christos G. Panayiotou
Department of Manufacturing Engineering
Boston University, Boston, MA 02215
cgc@bu.edu, gsun@bu.edu, panayiot@bu.edu

Abstract

We use Stochastic Fluid Models (SFM) for control and optimization of communication networks in which detailed discrete event models become impractical. By analyzing a SFM for a threshold-based admission control problem on a single node, we derive gradient estimators which can subsequently be used on the actual system. It is shown that these gradient estimators are unbiased and independent of all traffic and service processes involved, including the traffic and service rates. Moreover, they enable us to develop simple optimization schemes that recover the optimal thresholds of the actual discrete event model.

1 Introduction

A natural modeling framework for communication networks is provided through discrete event systems; in particular, queueing systems. However, the huge traffic volume that networks are supporting today makes such models highly impractical. Intuitively, one expects that the traffic volume can be exploited to model a packet flow as “fluid” and this has led to the development of Stochastic Fluid Models (SFM) for high-speed networks with bursty traffic, e.g., [1].

For the purpose of performance analysis with Quality of Service (QoS) requirements, the accuracy of SFMs depends on traffic conditions, the structure of the underlying system, and the nature of the performance metrics of interest. In this paper, our goal is to explore the use of SFMs for the purpose of *control and optimization* rather than *performance analysis*. In this case, it is not unreasonable to expect that one can identify the solution of an optimization problem based on a model which captures only those features of the underlying

“real” system that are needed to lead to the right solution, but not necessarily estimate the corresponding optimal performance with accuracy. Even if the exact solution cannot be obtained by such “lower-resolution” models, one can still obtain near-optimal points that exhibit robustness properties with respect to certain aspects of the model they are based on. Such observations have been made in several contexts (e.g., [2]), including recent results related to SFMs reported in [3] where for various optimization problems in queueing systems a connection between the SFM and queueing-system-based solution is established.

With this in mind, we have restricted ourselves here to admission control problems based on threshold parameters, so that the optimization problem of interest involves the determination of a threshold that minimizes some cost function. In a real network, thresholds are integer-valued, whereas in a SFM of the same system they are treated as real variables. Our first finding is that solving an optimization problem based on a SFM yields solutions that are close to those obtained if one were to solve the same problem using the actual discrete event system model. This motivates the need to develop efficient SFM-based techniques for solving such problems. Since such techniques often rely on gradient information, estimating the gradient of a given cost function with respect to the aforementioned threshold parameters in a SFM is an essential task for which Perturbation Analysis (PA) methods [4] are suitable. As recent work in [5], [6], [7] has shown, it is possible to derive PA estimators in the context of SFMs and establish their unbiasedness. In [6], for example, it is shown that Infinitesimal Perturbation Analysis (IPA) yields remarkably simple sensitivity estimates for packet loss metrics with respect to buffer size parameters. A particularly attractive aspect of such PA estimators is that they are obtained on line based on directly observable data (or data obtained by simulating the system as a SFM). Moreover, as we will show, these estimators are often independent of the system parameters as well, a crucial feature when detailed network data collection is not easily accomplished or when these parameters tend

¹This work was supported in part by the National Science Foundation under Grants ACI-98-73339 and EEC-0088073, by AFOSR under contract F49620-01-0056, by the Air Force Research Laboratory under contract F30602-99-C-0057 and by EPRI/ARO under contract WO8333-03.

to vary over time. This feature further enables us to use the IPA estimators derived for a SFM, but with data obtained directly from the actual system. This combination of SFM-based gradient estimators and real system data allows us to obtain actual optimal thresholds.

The contributions of the present paper are the following. First, we motivate our approach with an admission control problem solved using a SFM (Section 2). Motivated by this fact, we develop IPA estimators for loss and queueing content metrics with respect to an admission control threshold parameter in a single node model (Section 3). We subsequently use these estimators to solve admission control problems and include several numerical examples (Section 3.2).

2 A Motivating Example

Consider a network node where admission control at the packet level takes place using a simple threshold-based policy (see Fig. 1): When a packet arrives and the queue length is below a given amount $K = 1, 2, \dots$, it is accepted; otherwise it is rejected. Let $\bar{L}(K)$ denote the expected loss rate, i.e., the expected rate of rejected packets at steady state, and $\bar{Q}(K)$ the mean queue length when the threshold is K . We then define the following cost function

$$J(K) = \bar{Q}(K) + R \cdot \bar{L}(K) \quad (1)$$

where R is a penalty associated with rejecting a packet. Thus, $J(K)$ captures the trade-off between providing satisfactory service (low mean delay reflected by low $\bar{Q}(K)$) and rejecting too many packets (high loss rate $\bar{L}(K)$). The packet arrival process is modeled as an ON-OFF source such that packets arrive at a peak rate α during an ON period, followed by no packet arrivals during an OFF period. The packet processing rate is β . For the example used here and illustrated in Fig. 2, the number of arrivals in each ON period is geometrically distributed with parameter $p = 0.05$ and arrival rate $\alpha = 1$; the OFF period is exponentially distributed with parameter $\mu = 0.05$; and the service rate is $\beta = 0.55$. Thus, the traffic intensity of the system is $\alpha(\frac{1}{\alpha p})/\beta(\frac{1}{\alpha p} + \frac{1}{\mu}) = 0.909$, where $\frac{1}{\alpha p}$ is the average length of an ON period and $\frac{1}{\mu}$ is the average length of an OFF period. By exhaustively simulating this queueing system and averaging over 25 sample paths with 50,000 time units for each sample path and estimating $J(K)$ over different values of K , we obtained the curve labeled ‘DES’ in Fig. 2, using a rejection penalty $R = 30$. One can see that the optimal threshold value in this example is $K^* = 6$.

Next, we adopt a simple SFM for the same system: treating packets as “fluid”, during an ON period the queue content $X(t)$ increases at a rate $\alpha - \beta$ (assumed

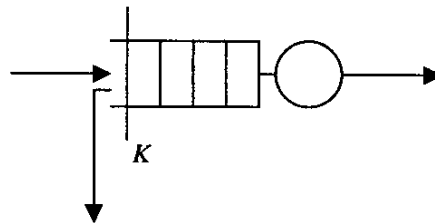


Figure 1: Admission Control in a Single Node

positive), while during an OFF period it decreases at a rate β . The cost function in this model is

$$J(u) = \bar{Q}(u) + R \cdot \bar{L}(u) \quad (2)$$

where $u \in \mathbb{R}_+$ is the threshold used to reject all fluid when the queue content reaches the level u . The corresponding mean queue content and expected loss rate are denoted by $\bar{Q}(u)$ and $\bar{L}(u)$ respectively. Simulating this model under the same ON-OFF conditions as before over different values of u results in the curve labeled “SFM” in Fig. 2. The important observation is that the two optima are close, whereas the difference in the actual cost estimates can be substantial (especially for a lightly loaded system.)

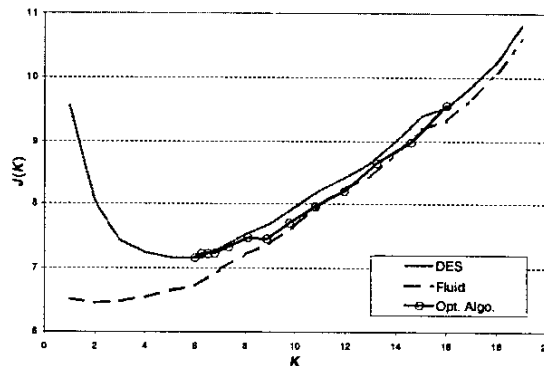


Figure 2: Cost v. Threshold Comparison for DES and SFM models

Based on this observation, we are motivated to study means for efficiently identifying solutions to problems formulated using a SFM. It is still difficult to obtain analytical solutions, however, since expressions for $\bar{Q}(u)$ and $\bar{L}(u)$ are unavailable, unless the arrival and service processes in the actual system are very simple. Therefore, one needs to resort to iterative methods such as stochastic approximation algorithms (e.g., [8]) which are driven by estimates of the gradient of a cost function with respect to the parameters of interest. In the case of the simple admission control problem above, we are interested in estimating dJ/du based on directly observed (or simulated) data. We can then seek to obtain u^* such that it minimizes $J(u)$ through an iterative

scheme of the form

$$u_{n+1} = u_n - \eta_n H_n(u_n, \omega_n^{SFM}), \quad n = 0, 1, \dots \quad (3)$$

where $H_n(u_n, \omega_n^{SFM})$ is an estimate of dJ/du evaluated at $u = u_n$ and based on information obtained from a sample path of the SFM denoted by ω_n^{SFM} . However, as we will see next, the simple form of $H_n(u_n, \omega_n^{SFM})$ that we will derive also enables us to apply it to ω_n^{DES} , a sample path of the real discrete event system:

$$K_{n+1} = K_n - \eta_n H_n(K_n, \omega_n^{DES}), \quad n = 0, 1, \dots \quad (4)$$

where K_n is the threshold used for the n th iteration. In other words, analyzing the SFM provides us with the structure of a gradient estimator whose actual value can be obtained based on data from the actual system. In Fig. 2, the curve labeled "Opt.Algo." corresponds to this process and illustrates how one can then recover the optimal threshold $K^* = 6$.

3 Admission Control for a Single Node

We begin with a model of a single node as a queueing system. In order to capture the bursty nature of incoming traffic, let $\alpha(t)$ denote the instantaneous input traffic rate at time t . The value of $\alpha(t)$ changes at time instants A_i , $i = 1, 2, \dots$, taking values from a finite set \mathcal{A}_α . For simplicity, we adopt an ON-OFF source model, so that $\mathcal{A}_\alpha = \{\alpha, 0\}$. i.e., the arrival rate is α during an ON period and 0 otherwise; it will become clear that our analysis readily extends to cases where \mathcal{A}_α contains more and arbitrary values.

The SFM for the same system can be represented by the same two parameters α and β , where α is the rate of a fluid entering the system during an ON period and β is the rate of the fluid exiting the system when its content is strictly positive. The stochastic component in this model lies in the duration of the ON and OFF periods, i.e., $\{A_i\}$, $i = 1, 2, \dots$, is a stochastic process describing the switches of the input rate between α and 0. It is also convenient to let α_i be the input rate over the interval $[A_i, A_{i+1}]$; clearly, $\alpha_i = \alpha$ if this is an ON period, otherwise $\alpha_i = 0$. We do not impose any restrictions on the interarrival time distribution during the ON period (other than $\alpha < \infty$) or on the distributions of the ON and OFF period durations. We also do not impose any restrictions on the service time distribution and let β denote the service rate. We assume that $\alpha > \beta$, otherwise a SFM model is trivial and buffer overflow phenomena are of limited interest.

We shall consider sample paths of the SFM whose length is defined by some final time $T < \infty$. Let $N(T) = \max\{i \geq 0 : A_i \leq T\}$ be the number of switches in $[0, T]$ and assume that $N(T) < \infty$ w.p.1, hence $E[N(T)] < \infty$. We define $X(t)$ to be the queue

content of the system at time $t \in [0, T]$ and note that $0 \leq X(t) \leq u$, where $u \in U \subset \mathbb{R}_+$ is the threshold used over the given sample path. The loss rate $\rho(t)$ is given by

$$\rho(t) = \begin{cases} \alpha - \beta & \text{if } X(t) = u \text{ and } \alpha > \beta \\ 0 & \text{otherwise} \end{cases}$$

For our purposes, we will concentrate on the process $\{X_i\}$, $i = 1, 2, \dots$, such that $X_i = X(A_i)$. i.e., the process defined by the queue content at the switching time instants A_i , $i = 1, 2, \dots$

Let $\mathcal{L}(u) : U \rightarrow \mathbb{R}$ be a sample function defined over a probability space (Ω, \mathcal{F}, P) and $J(u) = E[\mathcal{L}(u)]$. Strictly speaking, we write $\mathcal{L}(u, \omega)$ to indicate that the sample function depends on the sample point $\omega \in \Omega$, but will drop ω unless it is necessary to stress this fact. The optimization problems we consider involve the determination of u^* that minimizes $J(u)$. In order to accomplish this task, we rely on estimates of $dJ(u)/du$ provided by the sample derivatives $d\mathcal{L}(u)/du$.

For the admission control problem of interest here, we define two sample functions, the *Loss* $L_T(u)$ and the *Work* $Q_T(u)$ as follows:

$$L_T(u) = \int_0^T \rho(t) dt \quad \text{and} \quad Q_T(u) = \int_0^T X(t) dt$$

based on which the cost function of interest is

$$J_T(u) = \frac{1}{T} E[Q_T(u)] + \frac{R}{T} E[L_T(u)]$$

In order to estimate the derivative $dJ_T(u)/du$ and use this estimate in (3), we seek the sample derivatives $dL_T(u)/du$ and $dQ_T(u)/du$.

The starting point in PA is to consider a *nominal* sample path under some threshold u and a *perturbed* sample path resulting from perturbing u by Δu , keeping the realization of the process $\{A_i\}$, $i = 1, 2, \dots$, unchanged. For simplicity, we limit ourselves to the case where $\Delta u > 0$. We then define

$$\Delta X_i(u, \Delta u) = X_i(u + \Delta u) - X_i(u), \quad (5)$$

where $X_i(u)$ denotes the nominal sample path queue content at time A_i and $X_i(u + \Delta u)$ denotes the perturbed sample path queue content at the same time. Similarly, we define perturbations for some additional sample path quantities as follows. Setting $A_0 = 0$, let

$$L_i(u) = \int_{A_{i-1}}^{A_i} \rho(t) dt \quad \text{and} \quad Q_i(u) = \int_{A_{i-1}}^{A_i} X(t) dt$$

be the loss and work respectively measured over an interval $[A_{i-1}, A_i]$ which corresponds to either an ON or an OFF period of the input process. Then,

$\Delta Q_i(u, \Delta u) = Q_i(u + \Delta u) - Q_i(u)$, $\Delta L_i(u, \Delta u) = L_i(u + \Delta u) - L_i(u)$. In addition, set

$$Y_{i+1}(u) = X_i(u) + (\alpha_i - \beta)[A_{i+1} - A_i] \quad (6)$$

and note that $(\alpha_i - \beta)[A_{i+1} - A_i]$ is simply the amount of change in the queue content over the time interval $[A_{i+1} - A_i]$. Therefore, $Y_{i+1}(u)$ is the queue content obtained at time A_{i+1} if the queue were allowed to become $> u$ or < 0 . We then define $\Delta Y_i(u, \Delta u) = Y_i(u + \Delta u) - Y_i(u)$. Finally, note that in a SFM, as in a standard queueing model, a sample path consists of busy periods separated by idle periods. During a busy period the queue content is strictly positive and every busy period is followed by an idle period during which $X(t) = 0$. Let T_b be the time instant when the b th busy period of the nominal sample path ends. We define perturbations in the ending times of busy periods as $\Delta T_b(u, \Delta u) = T_b(u + \Delta u) - T_b(u)$, $b = 1, 2, \dots$. In what follows, we shall drop the arguments of all quantities ΔX_i , ΔY_i , ΔL_i , ΔQ_i , ΔT_b .

3.1 Perturbation Analysis

Let us consider a typical busy period indexed by $b = 1, 2, \dots$, and all possible events that can take place in it, so as to determine how associated perturbations are either generated (due to Δu) or propagated from the previous event. A busy period is initiated by an event that starts an ON period at some time A_i , and let us assume that $\Delta X_i = 0$ (we will discuss later the case where $\Delta X_i \neq 0$). The next event in the busy period is necessarily one that initiates an OFF period (we shall henceforth refer to this as an "OFF event") and there are two possible cases:

Case I: $Y_{i+1}(u) \leq u$. In this case, $Y_{i+1}(u)$ is given by (6) and we have: $X_{i+1}(u) = Y_{i+1}(u)$, $L_i(u) = 0$, $Q_{i+1}(u) = \frac{1}{2}Y_{i+1}(u)[A_{i+1} - A_i]$. Clearly, $\Delta X_{i+1} = \Delta Y_{i+1} = \Delta Q_{i+1} = \Delta L_{i+1} = 0$.

Case II: $Y_{i+1}(u) > u$. In this case, the queue content in the perturbed path can increase beyond u up to the perturbed threshold value $u + \Delta u$. Then,

$$\Delta X_{i+1} = \Delta u \quad (7)$$

$$\Delta L_{i+1} = -\Delta u \quad (8)$$

provided that Δu is such that $0 < \Delta u \leq Y_{i+1}(u) - u$; we will consider the case where $\Delta u > Y_{i+1}(u) - u$ below. The perturbation in work, ΔQ_{i+1} , is given by

$$\Delta Q_{i+1} = \Delta u \frac{Y_{i+1} - u}{\alpha - \beta} - \frac{(\Delta u)^2}{2(\alpha - \beta)}.$$

Let the length of the overflow interval in the nominal path be F_i and note that $F_i = (Y_{i+1} - u)/(\alpha - \beta)$. We therefore obtain

$$\Delta Q_{i+1} = F_i \Delta u + o(\Delta u) \quad (9)$$

where we use the notation $o(\Delta u) = O((\Delta u)^2)$. Assuming that an overflow interval has occurred in the observed busy period, (7)-(9) describe the perturbation generation process and how it affects the loss and work metrics of interest.

If $\Delta u > Y_{i+1}(u) - u = (\alpha - \beta)F_i$, then the perturbation reduces to $\frac{1}{2}(\alpha - \beta)F_i^2 < \frac{1}{2(\alpha - \beta)}(\Delta u)^2 = o(\Delta u)$ and we get

$$\Delta X_{i+1} = (\alpha - \beta)F_i \quad (10)$$

$$\Delta L_{i+1} = -(\alpha - \beta)F_i \quad (11)$$

$$\Delta Q_{i+1} = o(\Delta u) \quad (12)$$

Using the standard notation $[x]^+ = \max(x, 0)$, we can combine (7)-(9) with (10)-(12) to write

$$\Delta X_{i+1} = \Delta u - [\Delta u - (\alpha - \beta)F_i]^+ \quad (13)$$

$$\Delta L_{i+1} = -\Delta u + [\Delta u - (\alpha - \beta)F_i]^+ \quad (14)$$

$$\Delta Q_{i+1} = F_i \left(\Delta u - [\Delta u - (\alpha - \beta)F_i]^+ \right) + o(\Delta u) \quad (15)$$

To keep notation simple, let us set

$$\Delta u_i = \Delta u - [\Delta u - \Delta X_{i-1} - (\alpha - \beta)F_b]^+ \quad (16)$$

where F_b is the length of the most recent overflow interval in the b th busy period prior to A_i and $F_b = 0$ if no such interval exists (in which case $\Delta u_i = 0$, since no perturbation can be generated due to Δu). Note that $\Delta X_{i-1} = 0$ in **Case II** above, but, in general, $\Delta X_{i-1} \geq 0$.

Next, suppose this overflow interval ends with an OFF event at time A_i . The next event is one that causes a switch to an ON interval (we shall henceforth refer to this as an "ON event") at time A_{i+1} and there are three cases to consider:

Case 1: $Y_{i+1}(u) \geq 0$. In this case we have:

$$\Delta Y_{i+1} = \Delta X_{i+1} = \Delta X_i = \Delta u_i \quad (17)$$

$$\Delta L_{i+1} = 0 \quad (18)$$

and ΔQ_{i+1} is given by

$$\Delta Q_{i+1} = \Delta X_i [A_{i+1} - A_i] = \Delta u_i [A_{i+1} - A_i] \quad (19)$$

Case 2: $Y_{i+1}(u) < 0$ and $Y_{i+1}(u) + \Delta Y_{i+1} \leq 0$. In this case, the b th busy period ends and it is followed by an idle interval of length I_b , which in turn ends at time A_{i+1} . Clearly,

$$\Delta X_{i+1} = 0 \quad (20)$$

$$\Delta L_{i+1} = 0 \quad (21)$$

The change in queue content is given by

$$\Delta Q_{i+1} = \Delta u_i [A_{i+1} - A_i - I_b] + o(\Delta u) \quad (22)$$

Recalling that the end of the b th busy period is denoted by T_b , we have $T_b = A_{i+1} - I_b$, therefore the expression above becomes

$$\Delta Q_{i+1} = \Delta u_i [T_b - A_i] + o(\Delta u) \quad (23)$$

Case 3: $Y_{i+1}(u) < 0$ and $Y_{i+1}(u) + \Delta Y_{i+1} > 0$. This represents a situation where the idle period is eliminated in the perturbed path, i.e., $I_b < \Delta u_i / \beta$. Then, the queue content perturbation becomes

$$\Delta X_{i+1} = \Delta X_i - \beta I_b = \Delta u_i - \beta I_b \quad (24)$$

while no loss is involved in either path:

$$\Delta L_{i+1} = 0 \quad (25)$$

and ΔQ_{i+1} is given by:

$$\Delta Q_{i+1} = \Delta u_i [T_b - A_i] + \Delta u_i I_b - \frac{1}{2} \beta I_b^2$$

Since $I_b < \Delta u_i / \beta$, we can write

$$\begin{aligned} \Delta X_i I_b - \frac{1}{2} \beta I_b^2 &\leq \Delta u_i \frac{\Delta u_i}{\beta} + \frac{1}{2} \beta \left(\frac{\Delta u_i}{\beta} \right)^2 \\ &= \frac{3}{2\beta} (\Delta u_i)^2 = o(\Delta u) \end{aligned}$$

so that

$$\Delta Q_{i+1} = \Delta u_i [T_b - A_i] + o(\Delta u) \quad (26)$$

Next, let us assume that the busy period does not end after the ON event, i.e., **Case 1** applies. Then, the next event is an OFF event at time A_{i+1} and we can repeat the analysis of **Cases I-II**, except that now $\Delta X_i > 0$. There are two cases to consider, which are the generalizations of **Cases I-II** accounting for the possibility of some perturbation present before an OFF event occurs:

Case 4: $Y_{i+1}(u) \leq u$. This is identical to **Case 1** and we get (17)-(19). Note that whether $Y_{i+1}(u) + \Delta Y_{i+1} \leq u$ or $Y_{i+1}(u) + \Delta Y_{i+1} > u$ is irrelevant.

Case 5: $Y_{i+1}(u) > u$. $\Delta X_{i+1} = \Delta u$ as in **Case II**. Once again, however, it is possible that $\Delta u > (\alpha - \beta)F_i + \Delta X_i$, so that we write, similar to **Case II**,

$$\Delta X_{i+1} = \Delta u_{i+1} \quad (27)$$

In addition, if $\Delta u > (\alpha - \beta)F_i + \Delta X_i$, then $\Delta L_{i+1} = 0 - (Y_{i+1} - u) = -(\alpha - \beta)F_i$. Otherwise, $L_{i+1}(u + \Delta u) = Y_{i+1} + \Delta X_i - u - \Delta u$, and we get $\Delta L_{i+1} = \Delta X_i - \Delta u$. Thus,

$$\Delta L_{i+1} = -(\Delta X_i - \Delta u) + [\Delta u - \Delta X_i - (\alpha - \beta)F_i]^+ \quad (28)$$

Regarding ΔQ_{i+1} , we have $\Delta Q_{i+1} = \Delta u [A_{i+1} - A_i]$, which we can generalize to

$$\Delta Q_{i+1} = \Delta u_i [A_{i+1} - A_i] + o(\Delta u) \quad (29)$$

We can now collect the results of this analysis to characterize the cumulative loss and work perturbations over a busy period.

Lemma 3.1 For any busy period ending with an idle interval of length I_b following an OFF event at time A_i ,

$$\Delta X_{i+1} = [\Delta u_i - \beta I_b]^+ \quad (30)$$

Lemma 3.2 For all $i = 1, 2, \dots$,

$$0 \leq \Delta X_i \leq \Delta u \quad (31)$$

$$-\Delta u \leq \Delta L_i \leq 0 \quad (32)$$

The proofs of these Lemmas as well as of following results are omitted but can be found in [9]. Next, recall that the endpoints of busy periods are denoted by T_b , $b = 1, 2, \dots$. The perturbation in T_b can be easily obtained by noticing that

$$\Delta T_b = \frac{\Delta u}{\beta} \quad (33)$$

provided that $\Delta u \leq \beta I_b$. If $\Delta u > \beta I_b$, then the b th and $(b+1)$ th busy periods are merged, which implies that ΔT_b includes the entire length of the $(b+1)$ th busy period. Once again, in view of (16), to account for the fact that the b th busy period may contain an overflow interval of length F_i with $\Delta u > (\alpha - \beta)F_i + \Delta X_{i-1}$, Δu in (33) can be replaced by $\Delta u_i \leq \Delta u$.

In order to characterize the cumulative loss and work perturbations at the end of a busy period, let B_j denote the j th occurrence of an overflow observed on the nominal sample path. In the discrete event model, this corresponds to an instant when a packet is rejected because, upon arrival, it encounters a queue length given by the threshold parameter K . In the SFM, this corresponds to the fluid queue content exceeding the threshold level u . Given a busy period initiated at time A_i and ending at time T_b , we are interested in the presence of at least one such time instant in (A_i, T_b) . In particular, let

$$B_b^* = \begin{cases} \min\{B_j : B_j \in (A_i, T_b), j = 1, 2, \dots\} \\ \text{if } B_j \in (A_i, T_b) \text{ for some } j \\ 0 \text{ otherwise} \end{cases} \quad (34)$$

Making use of the standard indicator function $\mathbf{1}[B_b^* > 0] = 1$ if $B_b^* > 0$ and zero otherwise, we have the following result.

Lemma 3.3 Consider a busy period $[A_j, T_b]$ with $\Delta X_j = 0$ and let $A_m < T_b \leq A_{m+1}$. Assuming $\Delta u_i = \Delta u$ for all $i = j, \dots, m$, the cumulative loss and work perturbations at the end of this busy period are given by

$$\Delta L_b = -\Delta u \mathbf{1}[B_b^* > 0] \quad (35)$$

$$\Delta Q_b = [\Delta u [T_b - B_b^*] + o(\Delta u)] \mathbf{1}[B_b^* > 0] \quad (36)$$

We can now derive unbiased derivative estimates for the performance metrics of interest.

Theorem 3.1 *The derivatives of $E[L_T(u)]$ and $E[Q_T(u)]$ are given by*

$$\frac{dE[L_T(u)]}{du} = -E \left[\sum_{b=1}^M \mathbf{1}[B_b^* > 0] \right] \quad (37)$$

$$\frac{dE[Q_T(u)]}{du} = E \left[\sum_{b=1}^M [T_b - B_b^*] \mathbf{1}[B_b^* > 0] \right] \quad (38)$$

where M is the number of busy periods contained in $[0, T]$, including a possibly incomplete last busy period.

An immediate implication of this Theorem is that the following are unbiased estimators of $dE[L_T(u)]/du$ and $dE[Q_T(u)]/du$ respectively:

$$\left[\frac{dE[L_T(u)]}{du} \right]_{est} = - \sum_{b=1}^M \mathbf{1}[B_b^* > 0] \quad (39)$$

$$\left[\frac{dE[Q_T(u)]}{du} \right]_{est} = \sum_{b=1}^M [T_b - B_b^*] \mathbf{1}[B_b^* > 0] \quad (40)$$

Both estimators are very simple to implement: (39) is merely a counter of all busy periods in which at least one overflow takes place; (40) requires in addition a timer for all busy periods in which at least one overflow takes place, in order to record the time between the occurrence of the first overflow and the end of that busy period. Neither estimator requires any knowledge of the traffic or processing rates, nor does it depend on the nature of the random processes involved. Thus, it is easy to see that they can be used for more complex input processes that may involve switches between multiple traffic rates as opposed to only α and 0.

3.2 Optimal Admission Control Using IPA Estimators

Returning to (3) we can now use (39) and (40) as a gradient estimator

$$H_n(u, \omega_n^{SFM}) = \frac{1}{T} \sum_{b=1}^M (-1 + R[T_b - B_b^*]) \mathbf{1}[B_b^* > 0] \quad (41)$$

evaluated over a sample path ω_n^{SFM} of length T that contains M busy periods, following which a control update can be performed through (3).

The interesting observation here is that the same estimator may be used in (4) as follows: If a packet arrives and is rejected at some time B_b^* in the b th busy period of the actual system under a given threshold K_n , then this contributes $\mathbf{1}[B_b^* > 0] = 1$ and the exact same expression as in (41) can be used to update the threshold:

$$K_{n+1} = K_n - \eta_n H_n(K_n, \omega_n^{DES}), \quad n = 0, 1, \dots \quad (42)$$

The table below shows numerical results obtained by applying this scheme under different parameter settings leading to different traffic intensities ρ and comparing the optimal threshold K^* obtained through exhaustive simulation and the threshold K_{SFM}^* obtained through (42) using the SFM-based gradient estimator in (41). The simulation model is the same as the one described in Section 2.

α	p	μ	β	ρ	K^*	K_{SFM}^*
1	0.1	0.1	0.55	0.909	6	6.88
1	0.05	0.05	0.55	0.909	6	5.97
2	0.05	0.1	1.1	0.909	9	10.15
1	0.1	0.1	0.8	0.625	9	10.05
1	0.05	0.05	0.8	0.625	9	9.02
2	0.05	0.1	1.6	0.625	17	17.00

References

- [1] D. Anick, D. Mitra, and M. Sondhi, "Stochastic theory of a data-handling system with multiple sources," *The Bell System Technical Journal*, vol. 61, pp. 1871–1894, 1982.
- [2] B. Mohanty and C. Cassandras, "The effect of model uncertainty on some optimal routing problems," *Journal of Optimization Theory and Applications*, vol. 77, no. 2, pp. 257–290, 1993.
- [3] S. Meyn, "Sequencing and routing in multiclass queueing networks. Part I: Feedback regulation," in *2000 IEEE International Symposium on Information Theory*, pp. 4440–4445, 2000. To appear in *SIAM J. Control and Optimization*.
- [4] Y. C. Ho and X. Cao, *Perturbation Analysis of Discrete Event Dynamic Systems*. Dordrecht, Holland: Kluwer Academic Publishers, 1991.
- [5] Y. Liu and W. Gong, "Perturbation analysis for stochastic fluid queueing systems," in *Proc. 38th IEEE Conf. Dec. and Ctrl.*, pp. 4440–4445, 1999.
- [6] Y. Wardi and B. Melamed, "IPA gradient estimation for the loss volume in continuous flow models," in *Proc. of Hong Kong Intl. Workshop on New Directions of Control and Manufacturing*, pp. 30–33, November 1994.
- [7] Y. Wardi and B. Melamed, "Variational bounds and sensitivity analysis of traffic processes in continuous flow models," *J. of Discrete Event Dynamic Systems*, 2001. to appear.
- [8] H. Kushner and D. Clark, *Stochastic Approximation for Constrained and Unconstrained Systems*. Berlin, Germany: Springer-Verlag, 1978.
- [9] C. Cassandras, Y. Wardi, B. Melamed, G. Sun, and C. Panayiotou, "On-line gradient estimation for control and optimization of stochastic fluid models," *IEEE Trans. on Automatic Control*, 2001. Submitted.