

Perturbation Analysis and Optimization of Multiclass Multiobjective Stochastic Flow Models

Chen Yao and Christos G. Cassandras

Division of Systems Engineering
and Center for Information and Systems Engineering
Boston University
Brookline, MA 02446
cyao@bu.edu, cgc@bu.edu

Abstract—Stochastic Flow Models (SFMs) are stochastic hybrid systems that abstract the dynamics of complex discrete event systems involving the control of sharable resources. SFMs have been used to date to study systems with a single user class or some multiclass settings in which performance metrics are not class-dependent. In this paper, we develop a SFM framework for multiple classes and class-dependent performance objectives in which we can analyze new, occasionally counterintuitive, phenomena and give rise to a new type of “induced” events that capture delays in the SFM dynamics. In the case of two classes, we derive Infinitesimal Perturbation Analysis (IPA) estimators for their derivatives and use them as the basis for on-line optimization algorithms that apply to the underlying discrete event system (not the SFM). This allows us for the first time in the use of SFMs to contrast system-centric and user-centric objectives.

Keywords: Stochastic Flow Model, Perturbation Analysis, Stochastic Hybrid System, Discrete Event System

I. INTRODUCTION

The study of Discrete Event Systems (DES) is based on well-developed modeling frameworks in which the system dynamics are driven by the occurrence of different events defined over some given event set [2]. When event occurrence rates get extremely high, however, analysis becomes prohibitively complex; even well-designed discrete event simulations have impractically slow execution times. In this case, one seeks alternative models through which the system dynamics are *abstracted* to an appropriate level that retains essential features enabling effective and accurate control and optimization. This is often the case in systems where random phenomena play different roles at different time scales and typically gives rise to stochastic hybrid system models [3] in which some event-driven dynamics are retained to capture switches between different “modes” while the remaining dynamics are abstracted into differential equations describing the system state evolution within each such mode.

Fluid models are an example of this abstraction process applied to a large class of DES, and especially useful in analyzing communication networks with large traffic volumes [8]. While in most traditional fluid models the flow rates

involved are treated as fixed parameters, a *Stochastic Flow Model* (SFM), as introduced in [5], has the extra feature of treating flow rates as *stochastic processes*. With virtually no limitations imposed on the properties of such processes, a new approach for sensitivity analysis and optimization was recently proposed, based on Infinitesimal Perturbation Analysis (IPA). The essence of this approach is the on-line estimation of gradients (sensitivities) of certain performance measures with respect to various controllable parameters. These estimates may be incorporated in standard gradient-based algorithms to optimize parameter settings of the underlying DES. However, IPA estimates become biased (hence unreliable for control purposes) when dealing with aspects of queuing systems such as multiple user classes, blocking due to limited resource capacities, and various forms of feedback control. The emergence of SFMs has rekindled the interest in IPA because SFMs allow us to circumvent these limitations, yielding simple unbiased gradient estimates of useful metrics even in the presence of blocking and a variety of feedback control mechanisms [1],[10]. When it comes to multiple user classes, IPA has been applied to problems where flows are differentiated in terms of admission to a system, but once admitted all flows are treated alike [9]. IPA for SFMs that can differentiate flow classes in terms of service processes has been a challenge. More importantly, developing IPA estimates for gradients of class-dependent metrics has been elusive. Recently, Chen et al [6] have studied a multiclass SFM to analyze a dynamic priority call center. This model breaks new ground by differentiating among flow classes even after they enter the system; however, the analysis is very specific to the call center application and hard to extend to a general multiclass SFM model. In addition, the IPA analysis is limited to states but not general performance metrics, and unbiasedness for the estimators derived is not established.

In this paper, a general multiclass SFM model is developed. Each flow class, indexed by i , is associated with a threshold parameter θ_i based on which incoming flow is allowed into the system as long as $x_i(t) \leq \theta_i$, where $x_i(t)$ is the flow content of class i , which is thus differentiated from all other classes unlike earlier models such as in [9]. Moreover, each class is associated with its own performance metrics, such as workload, throughput, or loss rate due to

The authors' work is supported in part by NSF under Grants DMI-0330171 and EFRI-0735974, by AFOSR under grant FA9550-07-1-0361 and FA9550-09-1-0095, and by DOE under grant DE-FG52-06NA27490.

overflow. This is an important new element in the analysis of SFMs, allowing us to study the difference between user-centric and system-centric optimization, and place resource contention problems in a game framework. From a modeling and IPA standpoint, our approach introduces “induced events” in our SFM which can result in a (potentially infinite) event chain, a new phenomenon in the study of perturbation analysis, which allows us to understand some counterintuitive behavior observed in these multiclass environments.

II. MULTICLASS STOCHASTIC FLOW MODEL (SFM)

The systems we are interested in studying are those where multiple users are competing for service at a single sharable resource. Each user defines a “class” of tasks that are randomly generated, placed in a common queue, and processed on a FCFS basis across all such classes. This setting gives rise to a number of interesting problems where it is essential to distinguish between user-specific and system-centric performance metrics. We describe next a SFM abstraction of such a system, limiting ourselves to two user classes. It will become clear that our analysis directly applies to three or more classes at the expense of added notation that does not add to the basic ideas and results developed. Associated with a two-class SFM abstraction (see Fig. 1) are several real-valued and non-negative random processes which are all defined on a common probability space (Ω, \mathcal{F}, P) . The arrival flow processes $\{r_i(t)\}$, $i = 1, 2$, characterize the arrival rates of tasks at time t and the service flow processes $\{\beta_i(t)\}$, $i = 1, 2$, characterize their service rates. The total service capacity process is denoted by $\{C(t)\}$ and, clearly, $C(t) = \beta_1(t) + \beta_2(t)$. The process $\{x_i(t)\}$, $i = 1, 2$, defines the (real-valued) fluid content of class i in the system. In addition, a controllable parameter θ_i is associated with flow class i ; this is a threshold used in limiting the inflow or, alternatively, a buffer capacity assigned to class i : When $x_i(t) \geq \theta_i$, some of the class i incoming flow is dropped, giving rise to the overflow or loss process $\{l_i(t)\}$ and the input flow process $\{\alpha_i(t)\}$, $i = 1, 2$. Finally, the output flow processes $\{v_i(t)\}$, $i = 1, 2$, characterize departing flow rates. We are interested in the behavior of this SFM over a finite time interval $[0, T]$. Regarding the arrival and service processes, we will impose no restrictions on them as far as the probability laws that characterize them are concerned, but will make the following assumption:

Assumption 1: W.p. 1, the arrival $r_i(t) \geq 0$, $i = 1, 2$, and service capacity $C(t) \geq 0$ functions are piecewise constant in the interval $[0, T]$.

A. SFM Dynamics

We define a vector $\mathbf{x}(t) = (x_1(t), x_2(t))'$ where $x_i(t)$ is the class i queue content with the dynamics:

$$\frac{dx_i(t)}{dt^+} = \begin{cases} 0 & x_i(t) = 0 \text{ and } r_i(t) \leq \beta_i(t) \\ 0 & x_i(t) = \theta_i \text{ and } r_i(t) \geq \beta_i(t) \\ \alpha_i(t) - v_i(t) & \text{otherwise} \end{cases} \quad (1)$$

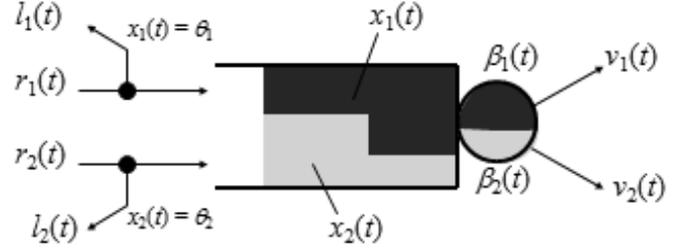


Fig. 1. A Two-Class Stochastic Flow Model (SFM)

where

$$\alpha_i(t) = \begin{cases} \beta_i(t) & x_i(t) = \theta_i \text{ and } r_i(t) \geq \beta_i(t) \\ r_i(t) & \text{otherwise} \end{cases} \quad (2)$$

$$v_i(t) = \begin{cases} r_i(t) & x_i(t) = 0 \text{ and } r_i(t) \leq \beta_i(t) \\ \beta_i(t) & \text{otherwise} \end{cases} \quad (3)$$

Thus, when $0 < x_i(t) < \theta_i$ we simply have $\frac{dx_i}{dt^+} = r_i(t) - \beta_i(t)$. When $x_i(t) = 0$ and $r_i(t) \leq \beta_i(t)$, the outflow rate is limited by the external arrival flow rate; similarly, when $x_i(t) = \theta_i$ and $r_i(t) \geq \beta_i(t)$, the inflow rate is limited by the service flow rate, leading to the loss rate, for $i = 1, 2$:

$$l_i(t) = r_i(t) - \alpha_i(t) = \begin{cases} r_i(t) - \beta_i(t) & x_i(t) = \theta_i \\ 0 & \text{and } r_i(t) \geq \beta_i(t) \\ & \text{otherwise} \end{cases} \quad (4)$$

We will use $x(t) = \sum_{i=1}^2 x_i(t)$ to denote the total system content at t . The crucial difference between a single class SFM, as in [4], and the two-class SFM is the behavior of the service rate $\beta_i(t)$. Whereas in the single-class model the service rate is independent of the system state, $\beta_i(t)$ in the two-class SFM depends on the queue contents and the inflow processes as explained next. Initially, the service rates are allocated proportional to the inflow rates, i.e.,

$$\beta_i(0) = C(0) \frac{\alpha_i(0)}{\sum_j \alpha_j(0)} \quad (5)$$

This allocation is maintained until there is a change in $\alpha_i(t) / \sum_j \alpha_j(t)$ at some time $t > 0$. When that happens, the total content $x(t)$ is the unprocessed workload under the initial service flow allocation. Let $\omega(t)$ denote the amount of time required to process this workload, at which point the new service rate allocation can take effect. Thus, the formal definition of $\omega(t)$ is through the relationship:

$$\int_t^{t+\omega(t)} C(\tau) d\tau = x(t) \quad x(t) > 0 \\ \omega(t) = 0 \quad x(t) = 0 \quad (6)$$

Finally, at time $t + \omega(t)$ the new allocation takes effect:

$$\beta_i(t + \omega(t)) = C(t + \omega(t)) \frac{\alpha_i(t)}{\sum_j \alpha_j(t)} \quad (7)$$

Therefore, in this SFM any event at t that causes a change in $\alpha_i(t) / \sum_j \alpha_j(t)$ is critical in that it “induces” another event

at $t + \omega(t)$ which results in a service rate allocation change. In the next section, we will formally define all events that can occur in the SFM. Here, we introduce the notation τ_k , $k = 1, 2, \dots$, to denote event occurrence times in increasing order of k . Using this notation, we define the set of *inflow change events* in the interval $[0, \tau_k]$:

$$F_k = \{m : \exists i \in \{1, 2\} \text{ s.t. } \alpha_i(\tau_m^-) \neq \alpha_i(\tau_m^+), m \leq k\} \quad (8)$$

To avoid degenerate cases where $C(\tau) = 0$ for all $\tau > t$, we will assume, whenever (6) is used, that $C(\tau) > 0$ for a sufficiently long time interval to ensure that $\omega(t) < \infty$.

The service flow allocation mechanism in (5)-(7) captures the FCFS nature of the underlying DES, as also noted in [6]. We can formally establish this fact by showing that if flow of both class types enters the SFM at time t , then it leaves the SFM at the same time, $t + \omega(t)$, regardless of the class type. This parallels the defining property of a FCFS policy, i.e., the waiting time of a customer arriving at t in a FCFS queue is the same regardless of its class. This result, stated as Lemma 2 below, rests on a monotonicity property expressed as Lemma 1. The proofs are fairly technical and are omitted, but they may be found in [11].

Lemma 1: If $t_1 \leq t_2$ and $x(t_1) > 0$, $x(t_2) > 0$, then $t_1 + \omega(t_1) \leq t_2 + \omega(t_2)$, where $\omega(t)$ was defined in (6).

Lemma 2: For $\omega(t)$ defined in (6),

$$\int_t^{t+\omega(t)} \beta_1(s) ds = x_1(t), \quad \int_t^{t+\omega(t)} \beta_2(s) ds = x_2(t)$$

Lemma 2 implies that any class i flow entering the SFM at t leaves at the same time $t + \omega(t)$ for $i = 1, 2$. If $r_1(t) = r_2(t) = 0$, $x(t) > 0$, then the lemma simply asserts that the time to deplete the current content of each class is the same, i.e., both class contents become zero simultaneously (unless of course $x_i(t) = 0$ for either $i = 1, 2$).

The presence of a delay $\omega(t)$ in the service flow allocation mechanism (5)-(7) also implies the need for additional state variables in the SFM dynamics. The role of these state variables is to provide “timers” triggered when an inflow change event occurs at t and then measure the amount of time until the queue content $x(t)$ is depleted. Thus, we define state variables $y_m(t)$, $m = 1, 2, \dots$, associated with events occurring at times τ_m , $m = 1, 2, \dots$, as follows:

$$\frac{dy_m(t)}{dt} = \begin{cases} -C(t) & \tau_m \leq t < \tau_m + \omega(\tau_m), m \in F_m \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

$$y_m(\tau_m^+) = \begin{cases} x(\tau_m) & y_m(\tau_m^-) = 0, m \in F_m \\ 0 & \text{otherwise} \end{cases}$$

Clearly, these state variables are only used for inflow change events, so that $y_m(t) = 0$ unless $m \in F_m$. Intuitively, $y_m(t)$ decreases from $x(\tau_m)$ at the rate of the service capacity $C(t)$ until this queue content is depleted, at which time $y_m(\tau_m + \omega(\tau_m)) = 0$ and the associated induced event takes place.

Similar to prior work on SFMs (e.g., [1],[10]), the class i queue content can be either empty, full, or neither. Accordingly, an interval $[\tau_k, \tau_l]$, $k < l$, over which $x_i(t) = 0$ for all $t \in [\tau_k, \tau_l]$ corresponds to an *empty period* (EP) for this

class and an interval $[\tau_k, \tau_l]$, $k < l$, over which $x_i(t) = \theta_i$ for all $t \in [\tau_k, \tau_l]$ corresponds to a *full period* (FP). A *boundary period* (BP) is either an empty or a full period; a *nonboundary period* (NBP) is a supremal interval during which $0 < x_i(t) < \theta_i$.

B. Event Classification

There are three types of events that can occur in the SFM of Fig. 1, as classified next.

1. Exogenous events. An event is *exogenous* if its occurrence time τ_k is independent of the controllable vector θ . In the two-class SFM, exogenous events correspond to changes (jumps, by Assumption 1) in the arrival flow rates $r_i(t)$, $i = 1, 2$, or the service capacity $C(t)$.

2. Endogenous events. An event occurring at time τ_k is *endogenous* if there exists a continuously-differentiable function $g_k(x, \theta)$ (see also [10]), such that

$$g_k(x(\theta, \tau_k), \theta) = 0 \quad (10)$$

In our model, the functions of interest are the boundaries defined by $x_i(\theta, t) = \theta_i$ or $x_i(\theta, t) = 0$, $i = 1, 2$, so that $g_k(x(\theta, t), \theta) = x_i(\theta, t) - \theta_i$ or $g_k(x(\theta, t), \theta) = x_i(\theta, t)$.

3. Induced events. We will refer to these as ω -events because they are all related to the definition of $\omega(t)$ in (6). An event of this type occurring at time τ_k is “induced” by an inflow change event at some time $\tau_m < \tau_k$, that is, any event (exogenous, endogenous, or itself an ω -event) such that some $\alpha_i(t)$, $i = 1, 2$, changes value at $t = \tau_m$, $m \in F_m$ as defined in (8). Thus, an event at time τ_k is an ω -event if there exists an event at τ_m , $m \in F_m$, such that

$$\tau_k(\theta) = \tau_m + \omega(\tau_m) > \tau_m$$

$$\text{and } \int_{\tau_m}^{\tau_m + \omega(\tau_m)} C(\tau) d\tau = x(\tau_m) > 0$$

It should be clear that an ω -event occurs at time $\tau_m + \omega(\tau_m)$ when the workflow $x(\tau_m)$ present at the time the event was induced becomes depleted and a service flow reallocation must result. If, however, $x(\tau_m) = 0$, by (6) we get $\omega(\tau_m) = 0$ and the event has no further effect on the SFM. We stress once again that an inflow change event that induces an ω -event may be exogenous, endogenous, or itself an ω -event; in the latter case, a chain of ω -events may be induced. Details and a description of counter-intuitive behavior resulting from an ω -event chain can be found in [11].

III. SFM PERFORMANCE OPTIMIZATION

An optimization problem for our SFM is defined by viewing $\theta = (\theta_1, \theta_2)$ as a controllable parameter vector and seeking to optimize performance metrics of the form

$$J(\theta; x(0), T) = E[\mathcal{L}(\theta; x(0), T)]$$

where $\mathcal{L}(\theta; x(0), T)$ is a sample function of interest evaluated in the interval $[0, T]$ with initial conditions $x(0)$. In this paper, we shall limit ourselves to the class-dependent loss volumes, $L_i(\theta; x(0), T)$, and average workloads, $Q_i(\theta; x(0), T)$, $i = 1, 2$, which will be explicitly defined in the sequel.

Given that we do not wish to impose any limitations on the defining processes $\{r_i(t)\}$ and $\{C(t)\}$ (other than mild technical conditions), it is infeasible to obtain closed-form expressions for $J(\theta; x(0), T)$. Therefore, we resort to iterative methods such as stochastic approximation algorithms (e.g., [7]), i.e. we seek to obtain θ^* minimizing $J(\theta; x(0), T)$ through an iterative scheme of the form

$$\theta_{i,n+1} = \theta_{i,n} - \eta_n H_{i,n}(\theta_n; x(0), T, \omega_n), \quad n = 0, 1, \dots \quad (11)$$

where $H_{i,n}(\theta_{i,n}; x(0), T, \omega_n)$ is an estimate of $\partial J / \partial \theta_i$ evaluated at $\theta = (\theta_{1,n}, \theta_{2,n})$ and based on information obtained from a sample path denoted by ω_n . In our work, we use the IPA sample derivative $\partial \mathcal{L} / \partial \theta_i$ as an estimate of $\partial J / \partial \theta_i$. Let N_T be the total number of events observed in a sample path over $[0, T]$. The average workload of flow class i , $i = 1, 2$, is

$$Q_i(\theta) = \frac{1}{T} \int_0^T x_i(t, \theta) dt = \frac{1}{T} \sum_{k \in \Omega_i} \int_{\tau_{k-1}}^{\tau_k} x_i(t, \theta) dt \quad (12)$$

with $\tau_0 = 0$, and Ω_i is the set of all non-empty periods (NEPs) for class i , defined as

$$\Omega_i = \{k : x_i(t) > 0 \text{ for all } t \in [\tau_{k-1}, \tau_k], k = 1, \dots, N_T\}$$

The average loss rate of flow class i , $i = 1, 2$, is

$$L_i(\theta) = \frac{1}{T} \int_0^T l_i(t, \theta) dt = \frac{1}{T} \sum_{k \in \Psi_i} \int_{\tau_{k-1}}^{\tau_k} [r_i(t) - \alpha_i(t)] dt \quad (13)$$

where we have used the fact that $l_i(t) = r_i(t) - \alpha_i(t)$ and $l_i(t) \geq 0$ only in FPs of class i , with the definition:

$$\Psi_i = \{k : x_i(t) = \theta_i \text{ for all } t \in [\tau_{k-1}, \tau_k], k = 1, \dots, N_T\}$$

In a single-class SFM, the average workload is an increasing function of a scalar threshold θ , and the loss rate is a decreasing function. Thus, we seek to strike a balance between these two metrics by determining θ that minimizes $J(\theta) = \gamma E[Q(\theta)] + E[L(\theta)]$ with some $\gamma > 0$. In our two-class SFM, however, we have a performance function that also reflects differences between classes, such as

$$J(\theta) = \gamma_1 E[Q_1(\theta)] \cdot \theta_1 + \gamma_2 E[L_1(\theta)] + \gamma_3 E[L_2(\theta)] + \gamma_4 E[Q_2(\theta)] \cdot \theta_2 \quad (14)$$

In addition, each class (user) may solve its own optimization problem with a performance metric of the form $J_i(\theta) = \gamma E[Q_i(\theta)] + E[L_i(\theta)]$, in which case we face a non-cooperative game setting.

Regardless of the optimization problem we choose to address, the starting point is the availability of estimates of $\partial J / \partial \theta_i$ which we obtain through IPA. It is clear from (12)-(13) that obtaining sample derivatives of $Q_i(\theta)$ and $L_i(\theta)$ derivatives requires the sample derivatives of the states $x_i(t, \theta)$ and of the event times $\tau_k(\theta)$ where the explicit dependence on the parameter θ is included for emphasis.

IV. IPA ESTIMATION

To simplify notation in the sequel, we define the following for all state and event time sample derivatives:

$$x'_{i,j}(t) \equiv \frac{\partial x_i(t)}{\partial \theta_j}, \quad y'_{m,j}(t) \equiv \frac{\partial y_m(t)}{\partial \theta_j}, \quad \tau'_{k,j} \equiv \frac{\partial \tau_k}{\partial \theta_j} \quad (15)$$

for $i, j = 1, 2$, and $k, m = 1, 2, \dots$

In the following, we will first discuss how to obtain state derivatives $x'_{i,j}(t)$, and then we will consider event time perturbations for each event type. Finally combining these two results will give us the IPA derivative estimation process. The mathematical derivation details are omitted, but can be found in [11].

A. State Derivatives

Let us rewrite the flow dynamics (1) over an interval $[\tau_{k-1}, \tau_k)$ as $\frac{dx_i(t)}{dt} = f_{i,k}(t)$ where either $f_{i,k}(t) = 0$ or $f_{i,k}(t) = r_i(t) - \beta_i(t)$. By studying the dynamics of $x'_{i,j}(t)$ defined in (15) for all $t \in [\tau_{k-1}, \tau_k)$ (see also [10]), we can get

$$x'_{i,j}(t) = x'_{i,j}(\tau_k^+) \quad t \in [\tau_k, \tau_{k+1}) \quad (16)$$

that is, the IPA derivative $x'_{i,j}(t)$ remains fixed in between consecutive events. In addition, taking advantage of the continuity of the queue content $x_i(t)$ at event time τ_k , we have

$$x'_{i,j}(\tau_k^+) = x'_{i,j}(\tau_k^-) + [f_{i,k}(\tau_k^-) - f_{i,k+1}(\tau_k^+)] \tau'_{k,j} \quad (17)$$

Thus, the queue content derivatives are piecewise constant, with jumps according to (17) at event times. It therefore suffices to use (17) to track them on an event by event basis.

Along the same lines, for the state variables $y_m(t)$, we have

$$y'_{m,j}(t) = y'_{m,j}(\tau_k^+) \quad t \in [\tau_k, \tau_{k+1}), \quad m \in F_k$$

and

$$y'_{m,j}(\tau_k^+) = y'_{m,j}(\tau_k^-) + [C(\tau_k^+) - C(\tau_k^-)] \tau'_{k,j} \quad (18)$$

If $y_m(\tau_k^-) = 0$ and $m \in F_{k-1}$, then, by definition, an induced event occurs at τ_k , and $\tau_k = \tau_m + \omega(\tau_m)$. Recalling (9), $y_m(t) = 0$ thereafter, so we also reset:

$$y'_{m,j}(\tau_k^+) = 0 \text{ if } \tau_k = \tau_m + \omega(\tau_m) \quad (19)$$

Finally, if an inflow change event occurs at τ_k , recalling (9), we have

$$y'_{k,j}(\tau_k^+) = \sum_{i=1}^2 x'_{i,j}(\tau_k^-) + \left[\sum_{i=1}^2 f_{i,k}(\tau_k^-) + C(\tau_k^+) \right] \tau'_{k,j} \quad (20)$$

In summary, (17) and (18)-(20) fully describe the propagation of the state derivatives from one event to the next, provided we can also evaluate all event time derivatives $\tau'_{k,j}$, $k = 1, 2, \dots$ as described next.

B. Event Time Derivatives

Recalling the event classification in Section II-B, we consider event time perturbations for each event type.

1. Exogenous events. By definition, all such events are independent of θ , therefore:

$$\tau'_{k,j} = 0, \quad j = 1, 2 \quad (21)$$

2. Endogenous events. In this case, (10) is in force and taking derivatives with respect to θ_j gives:

$$\frac{\partial g_k(\mathbf{x}, \tau_k, \theta)}{\partial \theta_j} + \frac{\partial g_k(\mathbf{x}, \tau_k, \theta)}{\partial \mathbf{x}} \left[\mathbf{x}'(\tau_{k,j}^-) + f_k(\tau_k^-) \cdot \tau'_{k,j} \right] = 0 \quad (22)$$

In the two-class SFM, we have either $g_k(\mathbf{x}(\theta, t), \theta) = x_i(\theta, t) - \theta_i$ or $g_k(\mathbf{x}(\theta, t), \theta) = x_i(\theta, t)$ for $i = 1, 2$. For convenience, we designate the corresponding endogenous events as follows: (i) A ρ_i event at τ_k is an event that initiates a FP for flow class i , i.e., $x_i(\tau_k^-) < \theta_i$, $x_i(\tau_k^+) = \theta_i$, and (ii) A σ event at τ_k is an event that initiates an EP i.e., $x_i(\tau_k^-) > 0$, $x_i(\tau_k^+) = 0$, $i = 1, 2$.

In the case of a ρ_i event at τ_k , (22) yields:

$$\tau'_{k,j} = \begin{cases} \frac{1 - x'_{i,j}(\tau_k^-)}{f_{i,k}(\tau_k^-)} & j = i \\ \frac{-x'_{i,j}(\tau_k^-)}{f_{i,k}(\tau_k^-)} & j \neq i \end{cases} \quad (23)$$

In the case of a σ event at τ_k , (22) yields:

$$\tau'_{k,j} = \frac{-x'_{i,j}(\tau_k^-)}{f_{i,k}(\tau_k^-)} \quad (24)$$

3. ω -events. Suppose an ω -event occurs at τ_k induced by an inflow change event at τ_m so that $\tau_k = \tau_m + \omega(\tau_m)$. Then, by definition, we must have $y_m(\tau_k^-) = 0$, and taking derivatives on both sides with respect to θ_j we get:

$$\tau'_{k,j} = \frac{y'_{m,j}(\tau_k^-)}{C(\tau_k^-)} \quad (25)$$

In summary, (21) for exogenous events, (23) for ρ_i events, (24) for σ events, and (25) for ω -events provide all the necessary event time derivatives, updated at all τ_k .

C. IPA Derivative Estimation Process

We can now combine our results to provide a complete description of the IPA derivative estimation process on an event by event basis. We proceed again using our event classification.

1. Exogenous events. Based on (21) in conjunction with (17) and (18)-(20):

$$\begin{aligned} \tau'_{k,j} &= 0 \\ x'_{i,j}(\tau_k^+) &= x'_{i,j}(\tau_k^-) \\ y'_{m,j}(\tau_k^+) &= y'_{m,j}(\tau_k^-) \quad m \in F_{k-1} \end{aligned} \quad (26)$$

2. Endogenous events. We consider the two types of endogenous events defined in the previous section.

2.1. ρ_i events, i.e., a FP for class i starts at time τ_k . In the case where $j = i$, (23) and (17) gives:

$$\tau'_{k,j} = \frac{1 - x'_{i,j}(\tau_k^-)}{r_i(\tau_k^-) - \beta_i(\tau_k^-)} \quad (27)$$

$$x'_{i,j}(\tau_k^+) = 1 \quad (28)$$

In the case where $j \neq i$, we get

$$\tau'_{k,j} = \frac{-x'_{i,j}(\tau_k^-)}{r_i(\tau_k^-) - \beta_i(\tau_k^-)} \quad (29)$$

$$x'_{i,j}(\tau_k^+) = x'_{i,j}(\tau_k^-) \quad (30)$$

Finally, using (18), we set

$$y'_m(\tau_k^+) = y'_m(\tau_k^-), \quad m \in F_{k-1}$$

2.2. σ events, i.e., an EP starts at time τ_k . Using (23) and (17), we get

$$\tau'_k = \frac{-x'_{i,j}(\tau_k^-)}{r_i(\tau_k^-) - \beta_i(\tau_k^-)} \quad (31)$$

$$x'_{i,j}(\tau_k^+) = 0 \quad (32)$$

In addition, we have

$$y'_m(\tau_k^+) = y'_m(\tau_k^-), \quad m \in F_{k-1}$$

3. ω -events. Suppose an ω -event occurs at τ_k induced by an inflow change event at τ_m so that $\tau_k = \tau_m + \omega(\tau_m)$. Using (25) gives

$$\tau'_{k,j} = \frac{y'_{m,j}(\tau_k^-)}{C(\tau_k^-)} \quad (33)$$

Therefore, (17) becomes

$$x'_{i,j}(\tau_k^+) = x'_{i,j}(\tau_k^-) + \quad (34)$$

$$\left[\alpha_i(\tau_k^-) - \alpha_i(\tau_k^+) + \beta_i(\tau_k^+) - \beta_i(\tau_k^-) \right] \frac{y'_{m,j}(\tau_k^-)}{C(\tau_k^-)}$$

In addition, we have

$$y'_m(\tau_k^+) = \begin{cases} 0 & \tau_k = \tau_m + \omega(\tau_m) \\ y'_m(\tau_k^-) & \text{otherwise} \end{cases}, \quad m \in F_{k-1}$$

Finally, if any event at τ_k is also an inflow change event, then we use (20) for the state variable y_k , where $\tau'_{k,j}$ is given by (26), (27), (29), (31), or (33) depending on the type of event that caused the inflow change.

Based on (26) through (34), we can evaluate all $x'_{i,j}(t)$ and $\tau'_{k,j}$ along a given sample path. We can then return to (12)-(13) and evaluate the performance metric derivatives $\frac{\partial Q_i(\theta)}{\partial \theta_j}$, $\frac{\partial L_i(\theta)}{\partial \theta_j}$ as described next.

Starting with $Q_i(\theta)$ in (12), we have for $i, j = 1, 2$:

$$\frac{\partial Q_i(\theta)}{\partial \theta_j} = \frac{1}{T} \sum_{k \in \Omega_i} x'_{i,j}(\tau_{k-1}) \cdot (\tau_k - \tau_{k-1}) \quad (35)$$

Similarly, for $L_i(\theta)$ in (13), we have:

$$\frac{\partial L_i(\theta)}{\partial \theta_j} = \frac{1}{T} \sum_{k \in \Psi_i} \left[[r_i(\tau_{k-1}^+) - \alpha_i(\tau_{k-1}^+)] (\tau'_{k,j} - \tau'_{k-1,j}) \right] \quad (36)$$

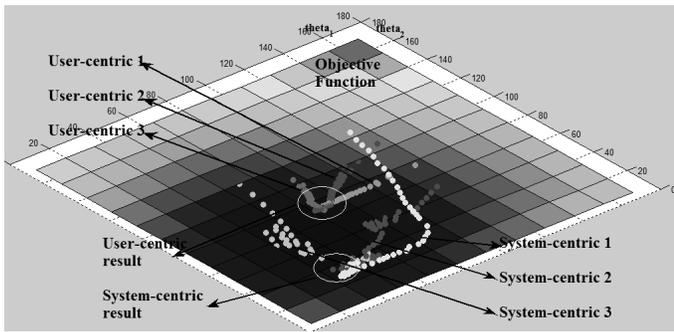


Fig. 2. Simulation result of system-centric and user-centric optimization

The unbiasedness of the IPA derivatives $\partial Q_i(\theta)/\partial\theta_j$ and $\partial L_i(\theta)/\partial\theta_j$ can be ensured by Assumption 1 and the following additional assumption.

Assumption 2. (a) For every $\theta \in \Theta$, w.p. 1, two events cannot occur at exactly the same time, unless one event induces the other, (b) W.p.1, no two processes $\{r_i(t)\}$, $\{\beta_i(t)\}$, $i = 1, 2$, have identical values during any open subinterval of $[0, T]$.

We can then establish the following result whose proof is omitted but may be found in [11].

Theorem 1. Under Assumptions 1-2, the IPA estimators $\partial Q_i(\theta)/\partial\theta_j$ and $\partial L_i(\theta)/\partial\theta_j$, $i, j = 1, 2$, are unbiased estimates of $dE[Q_i(\theta)]/d\theta_j$ and $dE[L_i(\theta)]/d\theta_j$, respectively.

V. DES OPTIMIZATION AND SIMULATION EXAMPLES

Recall that the SFM is used to derive performance sensitivity estimates that would otherwise not be possible for the actual DES. Let $J_T^{DES}(\theta)$ be some performance function of the DES and we seek to determine an optimal θ^* to minimize the $J_T^{DES}(\theta)$ through a standard stochastic approximation algorithm [7] as in (11):

$$\theta_{n+1} = \theta_n - \eta_n H_n(\theta_n, \omega_n^{DES}) \quad (37)$$

where $H_n(\theta_n, \omega_n^{DES})$ is an estimate of $dJ_T^{DES}(\theta)/d\theta$ at θ_n which is unavailable. Instead, we use the IPA estimator $dJ_T^{SFM}(\theta)/d\theta$ from the SFM. This requires identifying all events defined in the SFM with observable events in the real DES. A description of how this is accomplished can be found in [11].

Figure 2 shows an example of applying our IPA estimates and (37) in optimizing a two-class FCFS queueing system (not its SFM counterpart). The actual objective function shown compressed in two dimensions is obtained by exhaustive simulations of this DES over all (θ_1, θ_2) pairs, averaging over multiple sample paths. This gives (approximately) $\theta^* = (40, 40)$. The three trajectories labeled ‘‘System-centric’’ are all results of implementing (37) using gradient estimates obtained by applying the IPA algorithm on a single sample path with different starting points. We can see that each converges to a point sufficiently close to the ‘‘true’’ optimal, illustrating the effectiveness of our method.

As mentioned in the introduction, an interesting aspect of a multiclass system is that one can expect differences between a user-centric and a system-centric optimization approach. In system-centric optimization, we use (14) as the objective function. In the user-centric optimization, class 1 and class 2 take turns in optimizing their own performance metric $J_i(\theta) = \gamma_{1,i}E[Q_i(\theta_i)]\theta_i + \gamma_{2,i}E[L_i(\theta_i)]$. In this game, each class has no information on the other’s performance and has no control over the threshold of the other flow. Figure 2 shows the difference in these two optimization perspectives. The three trajectories labeled ‘‘User-centric’’ correspond to the same objective function as the system-centric approach but under user-centric optimization. One can see that the trajectories also converge to a common point, but the cost is larger than the point under the system-centric perspective. This gap actually implies the inefficiency of ‘‘selfish play’’ in such a resource-sharing system, sometimes referred to as ‘‘the price of anarchy’’.

VI. CONCLUSIONS AND FUTURE WORK

The two-class SFM we have studied opens up a spectrum of possibilities for studying systems with multiple user-specific objectives and reveals new perturbation dynamics not previously seen in SFMs. The IPA process remains relatively simple to implement, although it is no longer possible to sum it up in simple expressions as in single-class SFMs. We are currently studying extensions to a serial network of multiclass SFMs.

REFERENCES

- [1] C. G. Cassandras. Stochastic flow systems: Modeling and sensitivity analysis. In *Stochastic Hybrid Systems (C.G. Cassandras, and J. Lygeros, Eds)*, pages 139–167. Taylor and Francis, 2006.
- [2] C. G. Cassandras and S. Lafortune. *Introduction to Discrete Event Systems, Second Edition*. Springer, 2008.
- [3] C. G. Cassandras and J. Lygeros, editors. *Stochastic Hybrid Systems*. Taylor and Francis, 2006.
- [4] C. G. Cassandras, G. Sun, C. G. Panayiotou, and Y. Wardi. Perturbation analysis and control of two-class stochastic fluid models for communication networks. *IEEE Trans. on Automatic Control*, 48(5):770–782, 2003.
- [5] C. G. Cassandras, Y. Wardi, B. Melamed, G. Sun, and C. G. Panayiotou. Perturbation analysis for on-line control and optimization of stochastic fluid models. *IEEE Trans. on Automatic Control*, AC-47(8):1234–1248, 2002.
- [6] M. Chen, J-Q. Hu, and M. C. Fu. Perturbation analysis of a dynamic priority call center. *subm. to IEEE Trans. on Automatic Control*, 2008.
- [7] H. J. Kushner and G. G. Yin. *Stochastic Approximation Algorithms and Applications*. Springer-Verlag, New York, NY, 1997.
- [8] B. Liu, Y. Guo, J. Kurose, D. Towsley, and W. B. Gong. Fluid simulation of large scale networks: Issues and tradeoffs. In *Proc. of the Intl. Conf. on Parallel and Distributed Processing Techniques and Applications*, pages 2136–2142, June 1999.
- [9] G. Sun, C. G. Cassandras, and C. G. Panayiotou. Perturbation analysis of multiclass stochastic fluid models. *Journal of Discrete Event Dynamic Systems: Theory and Applications*, 14:267–307, 2004.
- [10] Y. Wardi, R. Adams, and B. Melamed. A unified approach to infinitesimal perturbation analysis in stochastic flow models: the single-stage case. *IEEE Trans. on Automatic Control*, 2009. To appear.
- [11] C. Yao and C. G. Cassandras. Perturbation analysis and optimization of multiclass multiobjective stochastic flow models, 2009. Technical Report, Div. of Systems Engineering, Boston University, <http://people.bu.edu/cyao/files/techreport.pdf>.