

# Efficient Simulation of Backlogs in Fluid Flow Lines

Dirk P. Kroese\*      Victor F. Nicola†

## Abstract

An Importance Sampling procedure to efficiently estimate overflow probabilities in continuous flow lines is described. The corresponding change of measure is found by generalizing the procedure that yields the change of measure for an ordinary fluid queue. Empirical results demonstrate the validity and effectiveness of the approach.

*Keywords:* Continuous flow line, production line, overflow probability, efficient simulation, importance sampling, likelihood ratio, optimal change of measure, eigenvalue equation, conjugate flow line, uniformization.

*AMS Subject Classifications (1991):* Primary 60K25; Secondary 90B22.

---

\*Department of Applied Mathematics, University of Twente, P.O. Box 217, 7500 AE Enschede, The Netherlands. Email: kroese@math.utwente.nl

†Department of Electrical Engineering, University of Twente, P.O. Box 217, 7500 AE Enschede, The Netherlands. Email: nicola@cs.utwente.nl

# 1 Introduction

A flow line is a tandem system of machines separated by buffers through which a stream of items/customers flows from one machine to the next. Flow lines are frequently encountered in manufacturing systems and other industrial processes, as well as in some computer and communication applications. A typical flow line is depicted in Figure 1.

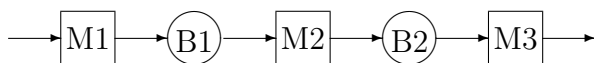


Figure 1: *A flow line with three machines and two buffers.*

The design of a flow line involves the dimensioning of its buffers in order to keep the probability of overflow sufficiently low. Therefore, the study of backlogs in such systems is of particular interest. Among the factors that determine the backlog behaviour are the processing speed and reliability of each machine, as well as the size of the associated buffer. Two performance measures are of particular interest: the probability of a buffer overflow during a “busy cycle”, and the stationary distribution of the content of a buffer. A comprehensive survey on flow lines may be found in [6]. Given the abundance of different flow lines, we remark that our model deals with a *continuous* flow line (i.e., the flow of products is modeled as a fluid flow) with *unreliable* machines (i.e., the machines are subject to service interruptions).

An exact analysis of flow lines is often not possible. Analytical results, mostly on stationary distributions, exist only for the most elementary systems. The 2-machine 1-buffer case (2-stage flow line) has been examined in [18], where the stationary distribution of the buffer content was found. For the 3-machine 2-buffer flow line with identical machines and (finite) buffers the joint stationary distribution of the buffers was found in [5]. Although in [9] more general 3-machine 2-buffer flow lines were considered, exact results were found only for a number of special cases which could be directly related to 2-stage flow lines. The fact that 3-stage flow lines are essentially more difficult to solve than 2-stage flow lines was demonstrated in [10], where a fluid tandem queue with on-off input was analyzed. This is basically a 3-stage flow line with one unreliable machine at the front of the line and two subsequent reliable machines. The joint stationary distribution of the content of the two

buffers was expressed in terms of integrals of modified Bessel functions.

An alternative approach to study flow lines is simulation. However, since backlog is typically a rare event, standard simulation is very inefficient, i.e., excessively long simulation time is required to achieve an acceptable relative accuracy. Importance sampling has been used successfully to speed up simulations involving rare events. However, this involves determining an appropriate change of measure to be used in the importance sampling simulations. In relatively simple queueing and fluid flow systems, such a change of measure may be obtained from an asymptotic large deviation analysis (see, for example, [4] and [13]). In [11], a large deviation analysis for the study of backlogs in a fluid flow line with unreliable machines is considered. The results will be used in the present paper to speed up simulations of backlogs in a flow line using an importance sampling procedure. Preliminary empirical results demonstrate the validity and effectiveness of our approach.

The remainder of this paper is organized as follows. In Section 2 we introduce the model. We will restrict ourselves to the most simple (3-stage) flow line for which no exact results exist. In Section 3 we briefly review the theory on fluid queues; in particular we recall the importance sampling procedure for such queueing systems. The main contribution of this paper is in Section 4, where we propose a method to efficiently simulate a 3-stage flow line, using the analogy between fluid queues and flow lines. In Section 5 we give a flow line example to show how the IS change of measure is derived in practice. In Section 6 we briefly discuss the uniformization approach and issues related to the implementation of importance sampling. And, in Section 7 we perform a simulation study of the flow line example using this change of measure. Finally, in Section 8 we give some conclusions and directions for further research.

## 2 The Model

We consider a continuous flow line with three machines and two buffers, as in Figure 1. Each machine is subject to failures and repairs. The life and repair times have exponential distributions and are all independent of each other. We denote the corresponding failure and repair rates by  $\lambda_i$  and  $\mu_i$ , respectively,  $i = 1, 2, 3$ . Each machine  $i$  has a specific *machine speed*  $\nu_i$ , which is the maximum rate at which products can be processed. Finally, the capacities of the two buffers are given by  $C_1$  and  $C_2$ , respectively. When a

buffer has reached its capacity any excess input stream of products is lost.

**Remark 1** Notice that the *actual* production speed of a machine  $i$  at a certain time may be lower than  $\nu_i$ . Consider for example the following situation. Assume  $\nu_1 > \nu_2 > \nu_3$ . Suppose machine 1 has failed and machines 2 and 3 are working. If the second buffer is not empty, machine 3 works at rate  $\nu_3$ . If the second buffer is empty but the first buffer is not, then machine 3 works at rate  $\nu_2$ . The machine does not process any products when both buffers are empty.

The state of the system at time  $t$  is determined by the state of each machine and the content of each buffer. Let  $M_i(t)$  be the state of machine  $i$  at time  $t$ ,  $i = 1, 2, 3$  (Working = 1, Failed = 0) and let  $Z_i(t)$  be the content of buffer  $i$  at time  $t$ ,  $i = 1, 2$ . Also define  $\mathbf{M}(t) = (M_1(t), M_2(t), M_3(t))$  and  $\mathbf{Z}(t) = (Z_1(t), Z_2(t))$ .

Obviously, the stochastic process  $(\mathbf{M}(t), \mathbf{Z}(t))$  is *regenerative*. A regeneration point is recognized when all buffers are empty, all machines are operational, and the input Markov modulating process (into the first buffer) enters a specified state. We assume that the process is *stable*, by which we mean that the expected length of a regeneration cycle is finite. In this case  $(\mathbf{M}(t), \mathbf{Z}(t))$  converges in distribution to some random vector  $(\mathbf{M}, \mathbf{Z})$ . The distribution of this vector is a major importance measure for the backlog behaviour of the system. However, our main concern will be the estimation of overflow probabilities of the *second*<sup>1</sup> buffer.

For a given buffer in the flow line, define a *regeneration overflow probability* as the probability that, starting from a regeneration point, the buffer content exceeds a given threshold before returning to level 0. We are primarily interested in the regeneration overflow probability of the second buffer. Large deviation results in [11] are used to derive an optimal change of measure, which we implement in an importance sampling procedure to efficiently estimate this rare event probability.

---

<sup>1</sup>The overflow behaviour of the first buffer is well known from the theory of 2-stage flow lines.

### 3 Efficient simulation of fluid queues

In this section we briefly review the theory on fluid queues. These systems are closely related to flow lines. In fact, we may view a flow line as a tandem system of fluid queues which are subject to service interruptions.

A fluid queue is a queueing system in which a reservoir is filled and depleted at rates which vary according to the current state of a regulating Markov chain. More precisely, let  $(X(t))$  be a continuous-time Markov chain with a finite state space, say  $E := \{1, \dots, m\}$ . This chain regulates the buffer content in such a way that net input rate is  $r_i$  whenever state  $i \in E$  is visited, provided that the buffer is non-empty. A fluid queue is therefore defined by two quantities: the generator (Q-matrix),  $Q = (q_{ij})$  say, of  $(X(t))$  and the set  $\{r_i\}$  of net input rates. It will be convenient to assemble the  $r_i$ 's in a diagonal matrix  $R := \text{Diag}(r_1, \dots, r_m)$ .

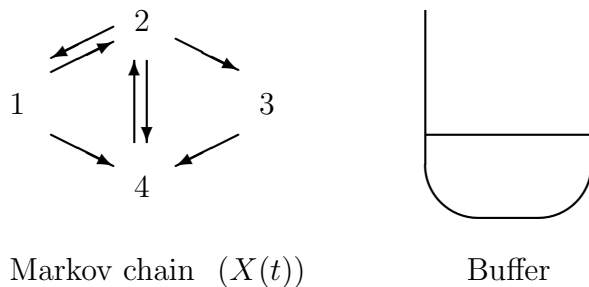


Figure 2: A fluid queue consists of a buffer and a regulating Markov chain.

In principle, the joint stationary distribution of the Markov chain and the content of the fluid reservoir can be determined analytically, see e.g., [1] and [16]. Overflow probabilities may be established in the same way. However, if the state space is large, then the numerical calculations needed to determine the overflow probabilities become rather involved<sup>2</sup> and simulation becomes a valid option. However, since overflow is typically a rare event, efficient simulation requires the use of an appropriate change of measure in an importance sampling procedure.

<sup>2</sup>Basically, all the eigenvalues of the matrix  $R^{-1}Q$  have to be calculated.

The following procedure establishes the optimal change of measure. Details can be found in Section 17 of [3], Chapter 3 of [13] or [15].

Firstly, determine the pair  $(\alpha, \mathbf{w})$ , satisfying the eigenvalue equation

$$Q \mathbf{w} = \alpha R \mathbf{w},$$

with  $\alpha$  the largest negative eigenvalue. Then, define a new Q-matrix  $Q^* = (q_{ij}^*)$  by putting

$$q_{ij}^* = q_{ij} \frac{w_j}{w_i}, \quad i \neq j,$$

where  $w_i$  is the  $i$ th element of  $\mathbf{w}$ .  $Q^*$  is called the *conjugate* Q-matrix of the fluid queue. The interpretation of  $Q^*$  is roughly the following: during an overflow period of the buffer, the driving Markov chain behaves like a chain with generator  $Q^*$  instead of  $Q$ . We now perform the simulation with  $Q^*$  instead of  $Q$ , and weigh the simulated events by the corresponding likelihood ratios. For example, if during a simulation run the driving chain has visited consecutively the states  $i_0, i_1, \dots, i_n$  with sojourn times  $s_0, s_1, \dots, s_n$ , then the likelihood ratio of such a run is given by

$$\frac{\prod_{k=0}^{n-1} q_{i_k i_{k+1}} e^{-\sum_{k=0}^n q_{i_k} s_k}}{\prod_{k=0}^{n-1} q_{i_k i_{k+1}}^* e^{-\sum_{k=0}^n q_{i_k}^* s_k}},$$

where  $q_i = \sum_{j \neq i} q_{ij}$  and  $q_i^* = \sum_{j \neq i} q_{ij}^*$ .

## 4 Efficient simulation of a flow line

In this section we establish the optimal change of measure that is to be used in the efficient simulation of the flow line of Figure 1. To simplify the analysis somewhat, we only consider flow lines with 3 machines and 2 buffers in which the last machine is *perfect* (cannot fail). Without loss of generality, we assume that  $\nu_1 > \nu_2 > \nu_3$  and the flow into Machine 1 is continuous at a fixed rate higher than  $\nu_1$ .

The key observation to make is that the second buffer may be viewed as the reservoir of a fluid queue that is regulated by the Markov process  $(X(t)) := (M_1(t), M_2(t), Z_1(t))$  (notice that  $M_3(t) = 1$ ). This driving process has a non-denumerable state space, so the theory of the previous section cannot be directly applied. However, we may reason analogously. For details we refer to [11].

The Markov process  $(X(t))$  is completely specified by its infinitesimal generator,  $\mathcal{Q}$  say.  $\mathcal{Q}$  is an operator which operates on vector-valued functions  $\mathbf{h} = (h_{00}, h_{01}, h_{10}, h_{11})^T$ , where the  $h_{ij} : \mathbb{R} \rightarrow \mathbb{R}$  are differentiable functions of  $x$ , the content of the first buffer. We have in particular,

$$(\mathcal{Q}\mathbf{h})(x) = Q\mathbf{h}(x) + \begin{cases} (\mathcal{B}_0\mathbf{h})(0), & x = 0, \\ (\mathcal{B}\mathbf{h})(x), & 0 < x < C_1, \\ (\mathcal{B}_{C_1}\mathbf{h})(C_1), & x = C_1, \end{cases} \quad (1)$$

where  $Q = (q_{ij,kl})$  is the Q-matrix<sup>3</sup> of  $(M_1(t), M_2(t))$ ,

$$\mathcal{B}_0 = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & \nu_1 \frac{d}{dx} & 0 \\ 0 & 0 & 0 & (\nu_1 - \nu_2) \frac{d}{dx} \end{pmatrix},$$

$$\mathcal{B} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -\nu_2 \frac{d}{dx} & 0 & 0 \\ 0 & 0 & \nu_1 \frac{d}{dx} & 0 \\ 0 & 0 & 0 & (\nu_1 - \nu_2) \frac{d}{dx} \end{pmatrix}$$

and

$$\mathcal{B}_{C_1} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & -\nu_2 \frac{d}{dx} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

The net input rate matrix  $R$  of the previous section is also replaced by an operator, namely  $\mathcal{R}$ , with

$$(\mathcal{R}\mathbf{h})(x) = \begin{cases} R_0 \mathbf{h}(0), & x = 0, \\ R \mathbf{h}(x), & x > 0, \end{cases}$$

where

$$\begin{aligned} R_0 &= \text{Diag}(-\nu_3, -\nu_3, -\nu_3, \nu_1 - \nu_3) \\ R &= \text{Diag}(-\nu_3, \nu_2 - \nu_3, -\nu_3, \nu_2 - \nu_3). \end{aligned}$$

---

<sup>3</sup>The order on the state space of  $(M_1(t), M_2(t))$  is lexicographic: 00,01,10,11.

By analogy with the previous section we determine the optimal change of measure of the flow line via the following procedure. (Empirical results support our claim of optimality, however, a formal proof is yet to be established.) Firstly, determine the pair  $(\alpha, \mathbf{w})$ , satisfying the eigenvalue equation

$$\mathcal{Q} \mathbf{w}(x) = \alpha \mathcal{R} \mathbf{w}(x), \quad 0 \leq x \leq C_1, \quad (2)$$

where  $\alpha$  is the largest negative eigenvalue. Next, define the *conjugate transition rates* as

$$q_{ij,kl}^*(x) = q_{ij,kl} \frac{w_{kl}(x)}{w_{ij}(x)}, \quad i, j, k, l \in \{0, 1\}, \quad (3)$$

where  $w_{ij}(x)$  is the  $ij$ -th element of  $\mathbf{w}(x)$ . Now perform the simulation with  $q_{ij,kl}^*(x)$  instead of  $q_{ij,kl}$ , and weigh the simulated events by the corresponding likelihood ratio.

Notice that the conjugate rates depend on  $x$ . The interpretation is that in the conjugate flow line (that is, during an overflow period of the second buffer) the failure and repair rates depend on the content of the first buffer.

Two questions remain. How to find the dominant eigenvalue  $\alpha$  and the corresponding function  $\mathbf{w}$ , and how to update the likelihood ratio. These issues are perhaps easiest to address in an example.

## 5 Flow line example

Consider the flow line of Figure 1 with the following parameters.  $\lambda_1 = 5$ ,  $\lambda_2 = 2$ ,  $\mu_1 = 1$ ,  $\mu_2 = 1$ ,  $\nu_1 = 3$ ,  $\nu_2 = 2$ ,  $\nu_3 = 1$ . The third machine is perfectly reliable. The buffer capacities are  $C_1 = 1$  and  $C_2 = \infty$ .

Notice that with any pair  $(\alpha, \mathbf{w})$  satisfying (2), the pair  $(\alpha, c\mathbf{w})$ , where  $c$  is a constant, satisfies the eigenvalue equation as well. We may therefore normalize  $\mathbf{w}$  such that  $w_{00}(0) = 1$ . Also, we define  $\kappa := w_{01}(0)$ .

For every  $\alpha$ , (2) represents a set of linear (differential) equations. The reader may verify that (2) is equivalent to the following:

$$\begin{aligned} w_{10}(0) &= 2 - \alpha - \kappa, \\ w_{11}(0) &= -2 + (3 - \alpha) \kappa. \\ w_{00}(x) &= \frac{w_{01}(x) + w_{10}(x)}{2 - \alpha}, \quad 0 \leq x \leq 1, \end{aligned}$$



and, for  $0 < x < 1$ ,

$$\begin{pmatrix} w_{01}(x) \\ w_{10}(x) \\ w_{11}(x) \end{pmatrix}' = \begin{pmatrix} -\frac{-4+\alpha+\alpha^2}{2(-2+\alpha)} & \frac{1}{2-\alpha} & \frac{1}{2} \\ \frac{5}{3(-2+\alpha)} & \frac{7-8\alpha+\alpha^2}{3(2-\alpha)} & -\frac{1}{3} \\ -5 & -2 & 7+\alpha \end{pmatrix} \begin{pmatrix} w_{01}(x) \\ w_{10}(x) \\ w_{11}(x) \end{pmatrix}, \quad (4)$$

and finally,

$$\begin{aligned} 5w_{00}(1) - 6w_{10}(1) + w_{11}(1) &= -\alpha w_{10}(1) \\ 5w_{01}(1) + 2w_{10}(1) - 7w_{11}(1) &= \alpha w_{11}(1). \end{aligned}$$

Thus,  $\mathbf{w}(1)$  is completely determined by  $\kappa$  and  $\alpha$ ; and these parameters can be found through the last two boundary conditions above.

We find<sup>4</sup>  $\alpha = -3.804641626$  and  $\kappa = 2.380068999$ . The  $3 \times 3$ -matrix in (4) has one real eigenvalue  $\gamma = 3.64901980591212$  and two complex conjugate eigenvalues  $\zeta \pm i\xi$ , with  $\zeta = 1.5510113201772$  and  $\xi = 1.11014777688726$ . It follows that the functions  $w_{ij}$  (defined in (2)) are of the form

$$w_{ij}(x) = a_{ij} e^{\gamma x} + e^{\zeta x} (b_{ij} \cos(\xi x) + c_{ij} \sin(\xi x)), \quad i, j \in \{0, 1\}.$$

The numerical values of the constants are given in Table 1.

$ij$	$a_{ij}$	$b_{ij}$	$c_{ij}$
00	0.00954550304762745	0.990454496952373	0.733972796830568
01	-0.0171811232870073	2.39725012189389	4.86407474791687
10	0.0725893476197914	3.35198327992163	-0.603625698973658
11	-0.130654877438759	14.3261714584643	4.38408359695598

Table 1: *The constants for  $w_{ij}(x)$ .*

The conjugate intensities now follow from (3), and are depicted in Figure 3. Notice that in the conjugate system the machine and repair rates depend on the state of the machines and the content of the first buffer.

---

<sup>4</sup>All constants in this section have been rounded off.

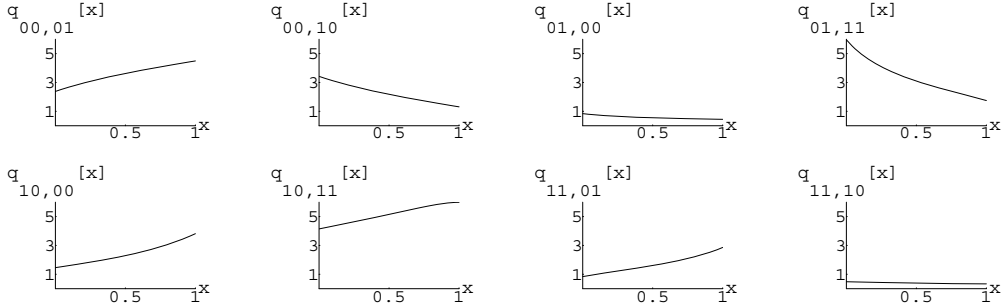


Figure 3: *Conjugate transition rates as functions of the buffer content,  $x$ .*

## 6 Uniformization and implementation issues

In Section 4 the optimal conjugate transition rates (to be used in importance sampling) are given as continuous functions of the content of the first buffer. That is, according to this change of measure, the failure and repair rates of machines are no longer constants but functions of the content of the first buffer. This suggests a suitable implementation of importance sampling using uniformization [14]. In this section we give a brief review of uniformization as a method for sampling from non-homogeneous Poisson processes. We also consider some implementation details in relation to our flow line problem.

Uniformization is a simple technique for sampling (i.e., simulating) the times of certain stochastic processes including non-homogeneous Poisson processes, renewal processes, or Markovian processes in continuous time on either discrete or continuous state spaces (see, for example, [7], [8], [12] and [17]). We describe it in the case of a non-homogeneous Poisson process ( $N(t)$ ,  $t \geq 0$ ) and assume that  $\lambda(t) \leq \beta$  for some finite constant  $\beta$ . Let  $T_n$  denote the time of the  $n$ -th event in a time homogeneous Poisson process ( $N_\beta(t)$ ,  $t \geq 0$ ) with a constant rate  $\beta$ . Then  $T_n$  is included (accepted) as an event in ( $N(t)$ ) with probability  $\lambda(T_n)/\beta$ , otherwise the point is not included (rejected). Rejected events are sometimes called pseudo events. Renewal processes can be simulated using uniformization as described above provided  $\lambda(t)$  is the hazard rate of the inter-event time distribution at time  $t$ . Uniformization can be generalized to cases in which the process being thinned is not a time homogeneous Poisson process (see [12]). For exam-

ple, at time  $T_{n-1}$ , we can let  $T_n = T_{n-1} + E_n$ , where  $E_n$  has an exponential distribution with rate  $\beta_n$ . The point  $T_n$  is then accepted with probability  $\lambda(T_n)/\beta_n$ . This requires only that  $\lambda(t) \leq \beta_n$  for all  $t \geq T_{n-1}$ .

As in Section 4, let us consider the flow line with 3 machines and two buffers in which the third machine is perfectly reliable. The transition rate matrix  $Q = (q_{ij,kl})$  is the generator matrix of  $(M_1(t), M_2(t))$ , where  $i, j, k, l \in \{0, 1\}$ . Let  $q_{ij} = \sum_{kl \neq ij} q_{ij,kl}$ , i.e.,  $q_{ij}$  is the total transition rate in state  $ij$ .

Uniformization can be used in flow line simulation as follows. Choose a constant uniformization rate  $\beta$ , such that  $\beta > \max_{i,j \in \{0,1\}} \{q_{ij}\}$ . Generate the times between uniformization events from an exponential distribution with a mean  $\beta^{-1}$ . In state  $ij$  ( $i, j \in \{0, 1\}$ ), a uniformization event is accepted as a real transition to state  $kl \neq ij$  with probability  $q_{ij,kl}/\beta$ . It follows that the uniformization event is rejected (i.e., it is a pseudo event) with probability  $(1 - \sum_{kl \neq ij} q_{ij,kl}/\beta) = (1 - q_{ij}/\beta)$ .

First, let us consider standard simulation to estimate the regeneration overflow probability  $p$  for the second buffer (as defined in Section 2). Generate  $n$  sample paths  $\omega_i, i = 1, 2, \dots, n$ , each starting with a pure regeneration (all machines are operational and all buffers are empty) and ending with the content of the second buffer hitting overflow (a rare event) or hitting zero (a typical event), whichever occurs first. For each sample path  $\omega_i$ , evaluate the indicator function  $I_i = I(\omega_i)$ , with  $I(\omega) = 1$ , if the sample path  $\omega$  includes a buffer overflow, otherwise,  $I(\omega) = 0$ . An unbiased estimator for  $p$  is given by  $\hat{p} = \frac{1}{n} \sum_{i=1}^n I_i$ . The relative error of this estimator is given by

$$\text{RE}(\hat{p}) = \frac{\sqrt{\text{VAR}(\hat{p})}}{E(\hat{p})} = \sqrt{\frac{1-p}{np}}.$$

It follows that  $\text{RE}(\hat{p}) \rightarrow \infty$  as  $p \rightarrow 0$ , which explains why standard simulation is not efficient for estimating very small values of  $p$ .

With importance sampling, the conjugate transition rates  $q_{ij,kl}^*(x)$  (as determined from (3)) are used instead of  $q_{ij,kl}$  to generate  $n$  samples of  $I(\omega)L(\omega)$ , where  $L(\omega)$  is the likelihood ratio associated with a sample path  $\omega$ . An unbiased estimator of  $p$  is given by  $\hat{p} = \frac{1}{n} \sum_{i=1}^n I_i L_i$ . For the above (optimal) change of measure, given by  $q_{ij,kl}^*(x), i, j, k, l \in \{0, 1\}$ , it can be shown that the relative error of this estimator remains bounded as  $p \rightarrow 0$ .

At the beginning of a sample path (pure regeneration), the likelihood ratio  $L$  is set to one, and importance sampling is turned on. The likelihood ratio is updated with each uniformization (real or pseudo) event in that sample

path as follows. In state  $ij$ , if an event is accepted as a transition to state  $kl \neq ij$ , then

$$L := L \times \frac{q_{ij,kl}}{q_{ij,kl}^*(x)}.$$

Otherwise, if the event is rejected (pseudo event), then

$$L := L \times \frac{(1 - q_{ij}/\beta)}{(1 - q_{ij}^*(x)/\beta)},$$

where  $q_{ij}$  and  $q_{ij}^*(x)$  are the total transition rates out of state  $ij$  in the original and the conjugate matrices,  $Q$  and  $Q^*$ , respectively. Importance sampling is turned off at the occurrence of either an overflow or an empty (second) buffer, whichever occurs first.

There is a considerable freedom in choosing the uniformization rate  $\beta$ , provided it is higher than the maximum total transition rate in the current state, according to both the original and the conjugate transition matrices,  $Q$  and  $Q^*$ , respectively. (Note that, with importance sampling, the total transition rate in a given state changes with the content of the first buffer.) An appropriate uniformization rate should be low enough to limit excessive generation of pseudo events, yet high enough to yield less noisy estimates of the likelihood ratio. Experiments have shown that a good choice of the uniformization rate is 5 to 10 times the maximum total transition rate in the current state. The uniformization rate can also be set independent of the state, provided it is sufficiently high; however, this may lead to excessive generation of pseudo events in states with low total transition rates.

## 7 Simulation results

In this section we perform a number of simulation experiments to demonstrate the validity of our approach. In particular, we use importance sampling to estimate the probability of overflow of the second buffer in the 3-stage flow line (described in Section 5) for several overflow levels,  $k$ . The optimal change of measure used in the IS method is a continuous function of the content of the first buffer,  $x$  (see (3) ), which follows directly from the constants in Table 1. Using an implementation based on uniformization (as described in the previous section) we compare standard simulation (SS) with importance sampling (IS) simulation. In each simulation run, we collect a fixed number of observations, each starting with a “pure” regeneration (all machines are

operational and all buffers are empty) and ending with the second buffer either reaching overflow or getting empty, whichever occurs first.

Let  $p$  be the probability that the content of the second buffer reaches the overflow level  $k$ , say, before hitting 0, starting with  $M_1(0) = M_2(0) = 1$  and  $Z_1(0) = Z_2(0) = 0$ .

In order to appropriately set the uniformization level  $\beta$ , a number of preliminary IS simulation runs are performed, in which the relative error  $\text{RE}(\hat{p})$  is determined for increasing values of  $\beta$  (see Table 2). Each simulation run is based on  $10^4$  observations. In Table 2, we see that initially the relative error  $\text{RE}(\hat{p})$  decreases sharply with an increase in  $\beta$ . Improvement in accuracy quickly levels off with further increase in  $\beta$ . A uniformization level  $\beta > 70$  yields a relative error close to the (asymptotic) minimum. However, an “optimal”  $\beta$  may be determined such that  $\beta \times (\text{RE}(\hat{p}))^2$  is minimal; this yields  $\beta = 80$  (i.e., about eight times the maximum total transition rate in any state of the original and conjugate processes.) A uniformization level higher than  $\beta = 80$  is not efficient, since it increases the simulation time without significantly improving the accuracy. Among other factors, the “optimal” uniformization rate depends on the (conjugate) hazard rate functions, which are used in importance sampling.

$\beta$	$\text{RE}(\hat{p}) (\times 10^{-2})$
10	27.6
20	5.29
30	3.81
40	3.31
50	2.65
60	2.75
70	2.35
80	2.14
90	2.08
100	2.07

Table 2: *Relative error  $\text{RE}(\hat{p})$  for increasing values of  $\beta$ .*

In Table 3 we list the estimates  $\hat{p}$  of  $p$ , using standard simulation and importance sampling, for various overflow levels  $k$ . For each estimate the half-width of the 95%-confidence interval is included, expressed as a percentage of

the estimate. All estimates are based on  $10^6$  observations, with a uniformization level  $\beta = 80$ . The results in Table 3 demonstrate a significant reduction of variance when using importance sampling, compared with standard simulation (for  $k = 4$ , no overflow events were observed with SS). Note that SS could be implemented more efficiently without uniformization; however, also in this case, experimental results show that IS (with uniformization) yields significant variance reduction for overflow levels as low as  $k = 2$ .

Notice also that the relative error increases slightly (at a linear rate) with the overflow level  $k$ . Numerical experiments (not reported here) suggest that such dependence does not occur when we use a sufficiently large  $\beta$ . Finally, it is interesting to note that

$$p \approx 0.7 e^{\alpha k},$$

for large  $k$ , where  $\alpha = -3.804641626$ , as determined in Section 5.

$k$	$\hat{p}$ (SS)	$\hat{p}$ (IS)
1	$1.53 \times 10^{-2} \pm 1.32\%$	$1.515 \times 10^{-2} \pm 0.32\%$
2	$3.58 \times 10^{-4} \pm 8.69\%$	$3.529 \times 10^{-4} \pm 0.37\%$
3	$6.00 \times 10^{-6} \pm 78.0\%$	$7.809 \times 10^{-6} \pm 0.43\%$
4	—	$1.739 \times 10^{-7} \pm 0.48\%$

Table 3: *Estimates of overflow probabilities using standard simulation (SS) and importance sampling (IS).*

## 8 Conclusions and further research

The backlog behaviour of a given buffer in a flow line with unreliable machines can be studied by viewing the buffer as a reservoir of a fluid queue driven by a Markov process. This modulating process depends on the state of the succeeding machine as well as all the preceding machine-buffer pairs. In general, this modulating process has a non-denumerable state space, with the number of continuous state variables equal to the number of the preceding buffers in the flow line.

Given a Markov-modulated characterization of its input, the backlog behaviour of a given buffer can be studied using large deviation results in an

importance sampling procedure. It is shown that the optimal change of measure is not fixed between discrete events (such as machine failure and repair events), but it must be varied continuously depending on the levels of all the preceding buffers. This suggests a suitable implementation of importance sampling using uniformization.

For a model with three machines and two buffers, empirical results demonstrate the validity and the effectiveness of the approach proposed in this paper. Further research is needed to investigate the feasibility of this and/or other approaches for flow lines consisting of more than two buffer-machine pairs. Also, an extension of the approach to non-Markovian failure/repair machine behaviour is of practical interest.

## References

- [1] D. Anick, D. Mitra and M.M. Sondhi, Stochastic theory of a data-handling system with multiple sources. *Bell System Tech. J.* **61** (1982) 1871–1894.
- [2] Asmussen, S. (1995) Stationary distributions for fluid flow models with or without Brownian noise. *Stochastic Models.* **11**(1) 21–49.
- [3] Asmussen, S. and Rubinstein, R.Y. (1995). Steady state rare events simulation in queueing models and its complexity properties. In *Advances in Queueing: theory, methods and open problems*, Dshalalow, J.H. (editor), CRC Press, New York, 429–461.
- [4] Chang, C.-S., Heidelberger, P., Juneja S., Shahabuddin, P. (1994). Effective bandwidth and fast simulation of ATMintree networks. *Performance Evaluation* **20**, 45–65.
- [5] Coillard, P. and Proth, J.-M. (1984). Sur l’effet de stocks tampons dans une fabrication en ligne. *Rev. Belge Statist. Inform. et Recherche Oper.* **24** 3–27.
- [6] Dallery, Y. and Gershwin, S. (1992). Manufacturing flow line systems: a review of models and analytical results. *Queueing systems* **12**, 3–94.
- [7] Gross, D. and Miller, D.R. (1984). The Randomization Technique as a Modeling Tool and Solution Procedure for Transient Markov Processes. *Operations Research* **32**, 343–361.

- [8] Jensen, A. (1953). Markov Chains as an Aid in the Study of Markov Processes. *Skand. Aktuarietidskr.* **36**, 87–91.
- [9] De Koster, M.B.M. and Wijngaard, J. (1986). A continuous flow model for three production units in series with buffers. *Operations Research proceedings DGOR* (Berlin: Springer-Verlag) 253–264.
- [10] Kroese, D.P. (1997). Analysis of a Markov-modulated fluid tandem queue. Memorandum No. 1367, University of Twente, Faculty of Applied Mathematics, Enschede, The Netherlands.
- [11] Kroese, D.P. (1997). How backlog builds up in a manufacturing flow line. In preparation.
- [12] Lewis, P.A.W. and Shedler, G.S. (1979). Simulation of Non-Homogeneous Poisson Processes by Thinning. *Naval Research Logistic Quarterly* **26**, 403–413.
- [13] Mandjes, M. (1996). Rare Event Analysis of Communication Networks, PhD Thesis, Free University, Amsterdam, The Netherlands.
- [14] Nicola, V.F., Heidelberger, P. and Shahabuddin, P. (1992). Uniformization and exponential transformation: Techniques for fast simulation of highly dependable non-Markovian systems. *Proc. of the Twenty-Third International Symposium on Fault-Tolerant Computing*. IEEE Computer Society Press, 130–139.
- [15] Ridder, A. (1996). Fast simulation of Markov fluid models. *J. Appl. Prob.* **33** 786–803.
- [16] L.C.G. Rogers, Fluid models in queueing theory and Wiener-Hopf factorization of Markov chains. *Ann. Appl. Probab.* **4** (1994) 390–413.
- [17] Shanthikumar, J.G. (1986). Uniformization and Hybrid Simulation/Analytic Models of Renewal Processes. *Operations Research* **34**, **4**, 573–580.
- [18] Zimmern, B. (1956). Etudes de la propagation des arrêts aléatoires dans les chaînes de production. *Rev. Statist. Appl.* **4** 85–104.