

# Large deviations without principle: Join the shortest queue

Ad Ridder<sup>1</sup>, Adam Shwartz<sup>2\*</sup>

<sup>1</sup> Vrije Universiteit, Amsterdam

<sup>2</sup> Electrical Engineering, Technion—Israel Institute of Technology

Received: date / Revised version: date

**Abstract** We develop a methodology for studying “large deviations type” questions. Our approach does not require that the large deviations principle holds, and is thus applicable to a large class of systems. We study a system of queues with exponential servers, which share an arrival stream. Arrivals are routed to the (weighted) shortest queue. It is not known whether the large deviations principle holds for this system. Using the tools developed here we derive large deviations type estimates for the most likely behavior, the most likely path to overflow and the probability of overflow. The analysis applies to any finite number of queues. We show via a counterexample that this system may exhibit unexpected behavior.

## 1 Introduction

The theory of large deviations is an important tool in the analysis of performance of computer communications networks. One encounters large deviations in probability problems at several levels, so in order to have a clear picture of the context in which we work, we briefly introduce the sample path large deviations considered in our study. For more extensive descriptions and the mathematical background we refer to [4–6, 10, 17, 19, 24].

Let  $y_n = (y_n(t) : t \in [0, T])$  for  $n = 0, 1, \dots$  be stochastic processes with sample paths in  $D([0, T], \mathbb{R}_+^K)$  (right continuous functions  $[0, T] \rightarrow \mathbb{R}_+^K$  with left limits; the space is equipped with the Skorohod metric).

---

\* Work of the first author was performed in part while visiting the Technion. Work of the second author was performed in part while visiting the Vrije Universiteit, Amsterdam, and was supported in part by Fund for the promotion of research at the Technion.

**Definition 1** We say that the sequence  $\{y_n\}$  satisfies the large deviations principle (LDP) with rate function  $I(\cdot)$  if  $I$  is nonnegative, lower semicontinuous, and for all  $x \in \mathbb{R}_+^K$  and for all closed  $C \subset D([0, T], \mathbb{R}_+^K)$ :

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_x (y_n \in C) \leq - \inf_{\phi \in C: \phi(0)=x} I(\phi),$$

and for all  $x \in \mathbb{R}_+^K$  and for all open  $G \subset D([0, T], \mathbb{R}_+^K)$ :

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_x (y_n \in G) \geq - \inf_{\phi \in G: \phi(0)=x} I(\phi).$$

(The subscript  $x$  means conditioning on  $y_n(0)$  so that  $\lim_{n \rightarrow \infty} y_n(0) = x$ .) We call  $\ell$  the local rate function if for any absolutely continuous  $\phi$ ,

$$I(\phi) = \int_{t_0}^{t_1} \ell(\phi(t), \dot{\phi}(t)) dt, \quad (1)$$

and  $I(\phi) = \infty$  if  $\phi$  is not absolutely continuous.

Let  $\{X(t) : t \geq 0\}$  be the Markov process describing the state of an exponential queueing model evolving in time, usually the number of jobs present at the various queues. Then one constructs a sequence of stochastic processes  $\{z_n\}$  by scaling time and space: the jump rates are speeded up by a factor of  $n$  whereas the jump sizes are diminished by a factor  $1/n$ , i.e.,  $z_n(t) = X(nt)/n$ . One studies these scaled queueing processes in the context of large deviations as introduced in definition 1. For instance for the simple M/M/1-queue, the LDP has been proved in [17, Section 11.4]. The rate function provides asymptotic or approximate expressions that can be used to develop algorithms and efficient simulation procedures for analysing and designing these systems. Therefore, much effort is devoted to the development of the theory of large deviations providing a useful expression for the rate function. However, this proves to be a difficult task in general. For example, the relatively broad study [5] does not cover queueing systems like the one we discuss here, while the general methods of identifying the rate function [1] do not cover any priority systems, and require checking rather complex hypotheses. At the other extreme, there are attempts to apply the theory to specific systems and questions. This is done by establishing first the validity of the large deviations principle for the specific system under investigation, and then calculating the probabilities of interest, see for example [2, 3, 8, 15, 23] and the applications in [17]. In other papers, the authors assume the large deviations principle to apply, and draw conclusions: see for example [11–13, 18]. For some cases, direct arguments allow to bypass the theory: see e.g. [20, 21] and the references in [9].

The queueing model of our interest is the “Join the shortest queue” (JSQ) system consisting of  $K$  queues each with a single exponential server. We allow for different rates at different servers. Arrivals from a single Poisson stream are routed to the “shortest weighted” queue, meaning the following.

A  $K$ -vector of positive weights  $(r_1, r_2, \dots, r_K)$  is given and when the state is the  $K$ -vector  $x = (x_1, x_2, \dots, x_K)$ ,  $x_i$  denoting the number of jobs present at queue  $i$ , an arrival is routed to the queue with the smallest value of  $x_i/r_i$ . The large deviations principle is not (yet) established for this simple queueing model because it does not fit in the framework mentioned above.

First we interpolate the scaled processes  $z_n$  between jump points and obtain (scaled) processes with continuous sample paths, because these are more convenient to work with. We do not distinguish between the piecewise-constant jump process and the piecewise-linear version, since the two are exponentially equivalent [17]. A path  $\phi$  is a continuous function  $[0, T] \rightarrow \mathbb{R}_+^K$ . Define the open  $\varepsilon$ -ball around a path  $\phi$  to be the collection of paths

$$B_\varepsilon(\phi) = \left\{ \psi : \sup_{0 \leq t \leq T} \|\psi(t) - \phi(t)\|_1 < \varepsilon \right\}.$$

Let  $x^0$  and  $x^T$  be two points in  $\mathbb{R}_+^K$  and  $\phi$  be a path with  $\phi(0) = x^0$ , and  $\phi(T) = x^T$ . We are concerned with the probability that  $z_n$  stays close to  $\phi$ , given that  $\lim_{n \rightarrow \infty} z_n(0) = x^0 = \phi(0)$ , denoted by  $\mathbb{P}_{x^0}(z_n \in B_\varepsilon(\phi))$ .

**Definition 2** We say that a path  $\phi$  is more likely than the path  $\psi$  if

$$\lim_{\varepsilon \downarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{x^0}(z_n \in B_\varepsilon(\psi)) \leq \lim_{\varepsilon \downarrow 0} \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{x^0}(z_n \in B_\varepsilon(\phi)).$$

We call  $\phi$  the optimal path from  $x^0$  to  $x^T$  if this holds for all  $\psi \neq \phi$ .

**Definition 3** We call  $I(\phi)$  the rate (or the cost) of path  $\phi$  if

$$I(\phi) := - \lim_{\varepsilon \downarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{x^0}(z_n \in B_\varepsilon(\psi)). \quad (2)$$

Events of practical interest may often be described in terms of sets of paths. For example,  $z_n(t)$  becomes large (i.e.  $|z_n(t)| \geq c$ ) if and only if the sample path  $z_n$  belongs to the set of paths with the same initial condition, and which satisfy  $|\phi(t)| \geq c$ . Thus our objective is to derive optimal (most likely) paths in the JSQ system for events of interest, show that their limits (2) exist, and compute their rates. Achieving this, we have established the asymptotic probability as well as most likely way an event occurs. Moreover, we have demonstrated that some basic techniques of large deviations can be useful even if it is not known whether the large deviations principle holds for the system of interest. Our major tools are coupling arguments, and large deviations for the Poisson process and the M/M/1 queue.

Previous work on the JSQ system—in relation to large deviations—covers the two-dimensional system only. For the symmetric system (where service rates at the queues are equal and both weights equal 1), the analysis in [17, Chapter 15.10] establishes the large deviations principle. However their analysis cannot be generalized to JSQ systems of higher dimensions or with asymmetric servers. A variation of the JSQ system which has been subject of large deviations studies, has a dedicated arrival stream for each

queue (besides the single stream that is routed to the smallest queue). This variant is analyzed in [1], where the large deviations principle is established. The analysis of this variant system does not allow the dedicated arrivals to have zero rate, thus the results of [1] cannot be translated to our JSQ system. Also [22] studies this variant JSQ system and computes the optimal path to overflow using coupling arguments, whereas [7, 14] derives a more precise description of the distribution on the way to overflow.

We have organized the paper as follows. In Section 2 we give a formal description of the model. We then restrict our attention to the two-dimensional system, so as to expose the ideas in a relatively simple setting. In Section 3 we describe the most likely behavior of this system, starting at any point. The way overflow occurs and the cost rates of overflow paths are given in Section 4. We show that, starting with an empty system, overflow always occurs by following the “weighted diagonal” at constant speed, possibly lingering at 0 before starting the excursion. For this case, we give an exact expression for the rate function. However, even in the two-dimensional case we show that, if the starting point is not the empty state, then it is possible that overflow occurs by emptying one queue (but not the other!) and then proceeding towards the “weighted diagonal.” We conclude in Section 5 with a description of the most likely behavior in any dimension, and comments on the overflow problem in higher dimensions.

## 2 The Join the Shortest Queue model

The two-dimensional Join the Shorter Queue (JSQ) system consists of two infinite queues, each with its own server. Service times are exponential, with parameter  $\mu_1$  and  $\mu_2$  respectively. The total service rate is  $\mu := \mu_1 + \mu_2$ . There are Poisson arrivals with rate  $\lambda$ . Let  $x_i$  denote the number of customers in queue  $i$ . An arrival is routed to one of the queues according to a control policy of the following type. The control is characterised by two positive numbers  $r_1, r_2$  such that whenever  $x_1/r_1 \leq x_2/r_2$ , the arriving customer is routed to queue 1, and otherwise to queue 2. To describe the control geometrically we imagine the line  $r \triangleq (r_1, r_2) \cdot \xi$ ,  $\xi \in \mathbb{R}_+$ . When the state  $x = (x_1, x_2)$  of the system is on or above the line  $r$ , arrivals join queue 1, and when the state is below the line, arrivals join queue 2. Let us call this line the control diagonal, or simply the diagonal. In accordance with the traditional JSQ terminology we say that arrivals join the shorter queue, although they actually join the “shortest weighted” queue. Generalizing to more than two queues, service rates are  $\mu_i$ ,  $i = 1, \dots, K$ , and the control is determined by the (strictly) positive numbers  $r_i$ ,  $i = 1, \dots, K$ . If the vector of queue sizes is  $x = (x_1, \dots, x_K)$  then an arrival is routed to the queue with the smallest value of  $x_i/r_i$ . We assume a fixed rule to break ties. We let  $i_x$  denote the queue to which an arrival will be sent, if the queue sizes are given by  $x$ . The particular choice is of no consequence to our analysis: we only assume that  $i_x = i_{\alpha x}$  for all  $\alpha > 0$ , that is, the choice depends

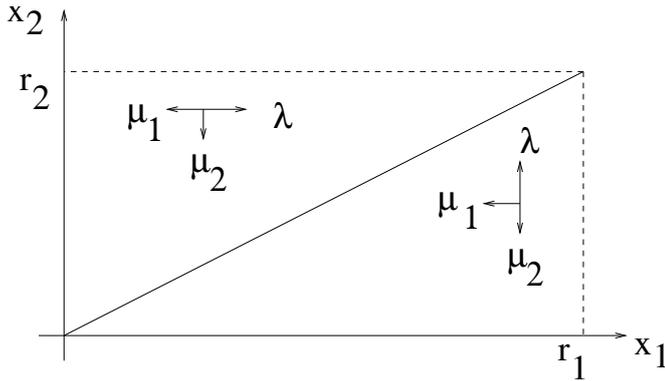


Fig. 1 Jump directions and rates in the JSQ system

only on the relative sizes of the queues. The case  $r_i = r_j$  corresponds to the traditional JSQ and then  $r$  is indeed the diagonal. Another special JSQ control is obtained by setting  $r_i = \mu_i$ : in this case the arriving customers join the queue with the smallest expected waiting time. Denote by  $e_i$  the unit vector in direction  $i = 1, \dots, K$ , that is,  $e_1 = (1, 0, \dots, 0)$  etc. We also denote  $e_{-i} = -e_i$ . Vectors are viewed as row vectors. We use the notation  $|x| = x_1 + \dots + x_K$ .

As introduced in Section 1 our scaled process is denoted by  $z_n$ : this is the process where all jumps are of size  $1/n$ , and all rates are multiplied by  $n$ . As before,  $z_{1,n} = z_n \cdot e_1$  is the first coordinate of the process  $z_n$ . Abusing notation we denote points of the scaled process again by  $x$ . They lie in the positive  $K$ -dimensional orthant  $\mathbb{R}_+^K$ . In our analysis we couple this process to several other processes, which we now define. Fix a point  $x$  and recall that  $i_x$  denotes the queue to which an arrival will be routed. The *local process*  $z_n^l$  is obtained by routing all arrivals to queue  $i_x$ , regardless of changes in queue sizes. The *reduced process*  $z_n^r$  is a scaled one-dimensional M/M/1 process with arrival rate  $n\lambda$ , service rate  $n\mu = n\mu_1 + \dots + n\mu_K$ , jumps of size  $1/n$  and starting point  $|x|$ . Note that the reduced process is stable (in the sense that it reaches 0 with probability one in finite expected time) if and only if  $\mu > \lambda$ . As we shall show, this is the stability condition for the JSQ system as well. Finally, the scaled *component process*  $z_n^c$  is the  $K + 1$ -dimensional process where coordinate  $i$  is a (scaled) Poisson process with rate  $n\mu_i$ , the  $K + 1$ st coordinate is a (scaled) Poisson process with rate  $n\lambda$  and all jumps are of size  $1/n$ . All four processes are coupled by using the same Poisson processes as the arrival and potential departure processes (potential meaning that when a jump occurs in the Poisson process while the associated queue is empty, no actual departure is triggered). Thus, every arrival entails a jump of size  $1/n$  in (a component of) all four processes, etc.

### 3 Most likely behavior

In this section we describe the most likely behavior of the process that starts at a point  $x$ . We provide a detailed analysis of the two-dimensional system. The generalization to  $K$  dimensions is given in Section 5. The behavior described below is the path that the process follows, with probability nearly equal 1, in the sense of Kurtz's Theorem [17, Theorem 5.3].

**Definition 4** Fix  $T > 0$ . We call  $z_\infty$  the most likely behavior for a sequence  $z_n$  of processes over  $[0, T]$  if the following hold. Given  $\varepsilon$  small enough, there is a constant  $C_1 > 0$  and a function  $C_2(\varepsilon) > 0$  so that

$$\mathbb{P}_x \left( \sup_{0 \leq t \leq T} |z_n(t) - z_\infty(t)| \geq \varepsilon \right) \leq C_1 e^{-nC_2(\varepsilon)} \quad \text{all } n > 0. \quad (3)$$

Note that the definition implies that, necessarily,  $z_\infty(0) = x$ .

#### 3.1 Most likely behavior in dimension 2

If  $x_1 > 0$ ,  $x_2 > 0$  and  $x$  is not on the ‘‘diagonal,’’ then at least locally one coordinate of the process is a Poisson process, and the other is the difference of two Poisson processes. Thus Kurtz's theorem applies [17], and we have

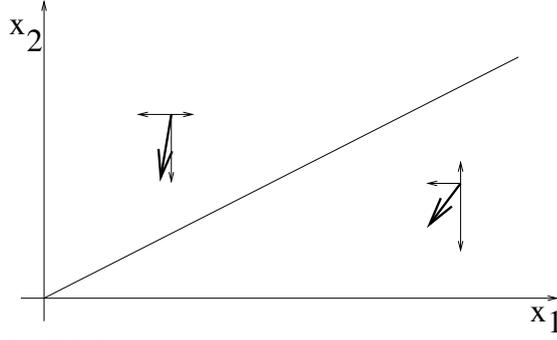
**Lemma 1** For  $x$  above  $r$  and off the boundary  $0 < x_1/r_1 < x_2/r_2$  define

$$z_\infty(t) = x + (\lambda - \mu_1)e_1 \cdot t - \mu_2 e_2 \cdot t. \quad (4)$$

Let  $T_x$  be the first time  $z_\infty(t)$  hits either the diagonal  $r$  or the vertical line  $(0, y)$ . Then the most likely path until time  $T_x$  is to follow  $z_\infty$ . Moreover,  $C_2$  is quadratic near  $\varepsilon = 0$ , i.e.,  $C_2 = O(\varepsilon^2)$  as  $\varepsilon \rightarrow 0$ . Finally, for any given  $\alpha > 0$ , the estimate (3) is uniform in  $\{x : T_x > \alpha\}$ . The analogous conclusion holds if  $x$  is below  $r$ .

The behavior depends both on the parameters and on the starting point. In general, if we start above  $r$  then, since arrivals all join queue 1,  $z_{2,\infty}$  can only decrease, so that  $T_x$  is finite. Hence, using elementary geometry we can establish all most likely paths with off-diagonal starting points. As an example, suppose that  $\mu_1 > \lambda$  and  $(\mu_2 - \lambda)/\mu_1 > r_2/r_1$ . Then the drifts above and below the diagonal are as plotted in Figure 2. So, when the starting point  $x$  is above  $r$  the most likely path moves down and closer to  $r$ ; when  $x$  is below  $r$  the most likely path moves down and away from  $r$ .

Consider diagonal starting points  $x \neq 0$ . The most likely path depends on the drifts above and below the diagonal and these are given by the parameters. The drift may be towards the diagonal in one part, and away in the other, as in Figure 2, or towards the diagonal in both parts. There is no combination of parameters so that both drifts are away from the diagonal. The two feasible possibilities are dealt with in the next two Lemmas.



**Fig. 2** Example of drifts

**Lemma 2** Suppose the starting point  $x \neq 0$  is on  $r$  and that the drift above the diagonal is towards and the drift below the diagonal is away from the diagonal (Figure 2). Let  $T_x \triangleq x_2(\mu_2 - \lambda)^{-1}$ . Fix  $T < T_x$  and define

$$z_\infty(t) = x - \mu_1 e_1 \cdot t + (\lambda - \mu_2) e_2 \cdot t.$$

Then  $z_\infty$  is the most likely path until time  $T$ ,  $C_2$  is quadratic near  $\varepsilon = 0$  and for any  $\alpha > 0$ , the estimate (3) is uniform in  $\{x : T_x > \alpha\}$ .

*Proof.* Let  $z_n^l$  denote the local process where arrivals join queue 2. If we couple the two queueing systems (that is, drive them with the same arrival and departure processes) then, at least until one process hits the boundary,  $z_{1,n}(t) \geq z_{1,n}^l(t)$  and  $z_{2,n}(t) \leq z_{2,n}^l(t)$ . This holds since some arrivals to queue 2 for the process  $z_n^l$  are sent to queue 1 in the JSQ system (when above  $r$ ). However,  $z_n^l$  satisfies the conditions of Kurtz's theorem and so (3) holds. This implies that the scaled JSQ process moves away from  $r$ . But once it does, its increments are exactly those of  $z_n^l$ , and another application of Kurtz's theorem establishes the Lemma. ■

**Lemma 3** Suppose the starting point  $x \neq 0$  is on  $r$  and that both drifts are towards the diagonal. Let

$$z_\infty^r(t) = x_1 + x_2 + (\lambda - \mu_1 - \mu_2) \cdot t.$$

If the system is stable, define  $T_x$  as the first time  $z_\infty^r(t) = 0$ . Otherwise, choose an arbitrary finite  $T_x$ . Fix  $T < T_x$  and define  $z_\infty$  through

$$z_{2,\infty}(t)/r_2 = z_{1,\infty}(t)/r_1 \quad \text{and} \quad z_{1,\infty}(t) + z_{2,\infty}(t) = z_\infty^r(t). \quad (5)$$

Then the most likely path until time  $T$  is  $z_\infty$ . Moreover,  $C_2$  is quadratic near  $\varepsilon = 0$ . Finally, for any  $\alpha > 0$ , the estimate (3) is uniform in  $\{x : T_x > \alpha\}$ .

*Proof.* The proof uses standard ideas which will be elaborated for the large deviations calculation, and so will only be described briefly here. First note that the two conditions in (5) define a unique path. By coupling the JSQ process with the reduced process we immediately have that the second condition must hold: that is, in the sense of Kurtz's theorem, the scaled process must satisfy this condition. Consider now an initial short period of time. Since the component process satisfies the conditions of Kurtz's theorem, we know that the process (with high probability) cannot move far. If it stays near  $r$ , we repeat the argument, to show that it stays near  $r$  until  $T$ . Once it moves away from  $r$ , it will be driven towards the diagonal as we have established earlier. Thus the process stays near  $r$  until  $T$  (in the sense of Kurtz theorem). Thus the first condition of (5) holds as well. ■

The boundary starting points remain. The following two Lemmas deal with it, first the  $x \neq 0$  case.

**Lemma 4** *Suppose  $x_1 = 0$  and  $x_2 > 0$ .*

- (i) *If  $\mu_1 \geq \lambda$  then (3) holds with  $z_\infty(t) = x - \mu_2 e_2 \cdot t$  for  $T < T_x = \frac{x_2}{\mu_2}$ .*
- (ii) *If  $\mu_1 < \lambda$  then (3) holds where  $z_\infty(t)$  is given by (4), until  $z_\infty$  meets  $r$ .*

The proof couples the JSQ process and the local process: see our report [16].

Note that we have now established our stability claim: the JSQ system is stable if  $\mu_1 + \mu_2 > \lambda$ .

**Lemma 5** *Suppose  $x = 0$ . If the JSQ is stable then the most likely behavior is to stay near 0. If it is unstable, then the most likely behavior is to follow  $r$ , so that  $z_{1,\infty}(t) + z_{2,\infty}(t) = (\lambda - \mu_1 - \mu_2) \cdot t$ .*

*Proof.* If the JSQ is stable then whenever it moves away from 0, the previous results show that it must move towards 0.

In the unstable case, the reduced system follows a straight line away from 0 with speed at least  $\lambda - \mu_1 - \mu_2 > 0$ . This is the speed if we are away from the boundaries. But as soon as we move away from 0, our previous results imply that the process moves towards  $r$ , and so away from the boundaries. Therefore the process moves along  $r$  with speed  $\lambda - \mu_1 - \mu_2$ . ■

#### 4 Most likely path to overflow: dimension 2

We return to our objectives stated in Section 1: finding optimal paths and their rates. We shall consider only optimal paths to overflow—to be defined shortly—since these are the most interesting ones. We provide an explicit formula for the case that the system starts empty. Let  $C = nc$  be some high level for queue sizes. Buffer overflow occurs when the system reaches any state where one or both queues exceed the high level:  $\{(x_1, x_2) : x_1 \geq nc \text{ or } x_2 \geq nc\}$ . Unless we start very close to the set of overflow states, this set is hit at the unique system state  $x = (r_1, r_2)n\beta$  where the diagonal crosses the boundary of the overflow set. For example,

if  $r_1 > r_2$  then  $\beta = c/r_1$ . In this section we assume explicitly that the JSQ system is stable, for otherwise the event of reaching overflow is not rare, and the most likely behavior brings the process to the mentioned overflow state (with probability equal nearly 1, as demonstrated in the previous section). Finally notice that  $(r_1, r_2)\beta$  is the overflow point for the scaled process. Summarizing, our objective is to analyse limiting properties of  $\mathbb{P}_{x^0}(z_n \in B_\varepsilon(\phi))$  for paths with  $\phi(0) = x^0$  and  $\phi(T) = (r_1, r_2)\beta$ .

Our first observation is that the first queue cannot grow below the diagonal and the second queue cannot grow above the diagonal, the only way to reach the overflow point  $(r_1, r_2)\beta$  is by reaching the diagonal and then following it in an upward direction (both coordinates increase). A path or a piece of path where both coordinates increase (strictly!) is called *increasing*. We can split our search of optimal path into two:

- Given a starting point  $x^0$  and a time  $T$ , how is the diagonal reached?
- Given a starting point on the diagonal and a time  $T$ , what is the optimal path that increases along the diagonal?

We deal with the second question: it will be clear that the same methods apply to the first. So, consider paths which start on the diagonal,  $r^0 := (r_1, r_2)\alpha$  for some  $0 < \alpha < \beta$ , increase along the diagonal, and end at the overflow point  $r^T := (r_1, r_2)\beta$  at time  $T$ . Notice that we restrict—for the moment—the starting point to avoid the point 0. The reason is that if the process stays near the diagonal starting at  $r^0$ , both servers are kept busy. This makes it easier to treat the process. The speed of increase during the time period  $[0, T]$  is allowed to be any positive function, but we claim that paths with constant speed are the most likely.

**Lemma 6** *Let  $\psi$  be a path from  $r^0$  to  $r^T$ , increasing along the diagonal. Let*

$$\phi(t) := r^0 + (r_1, r_2)(\beta - \alpha)t/T, \quad 0 \leq t \leq T.$$

*Then  $\phi$  is more likely than  $\psi$ .*

*Proof.* Recall  $|x| := x_1 + x_2$ . We show that the probability for the JSQ process to stay near  $\phi$  is (asymptotically) the same as the probability for the reduced process to stay near  $|\phi|$  and, in addition, it is larger than the probability to stay near any other path  $\psi$  that stays on the diagonal.

Recall that the coupled reduced process  $z_n^r$  is an M/M/1 queue with arrival rate  $\lambda$  and service rate  $\mu = \mu_1 + \mu_2$ , and recall that M/M/1 process satisfies the LDP [17, Section 11.4]. Paths with constant speed in a one-dimensional process are said to be straight lines. For  $\varepsilon < \min\{r_1^0, r_2^0\}$ , since  $\psi$  is increasing and  $|\psi(0)| = |r^0| \neq 0$ , we have that  $|\psi(t)| \geq \varepsilon > 0$  and

$$\{z_n \in B_\varepsilon(\psi)\} \subset \{z_n^r \in B_\varepsilon(|\psi|)\} \subset \{z_n^r > 0\}. \quad (6)$$

The inclusions (6) hold since if the JSQ process stays near  $\psi$  then the coupled reduced process satisfies these conditions. However, for the reduced process we have an LDP with a convex local rate function inducing that

straight lines are more likely than other increasing paths [17, Lemma 5.16]. Here there is only one straight line from  $|r^0|$  to  $|r^T|$ , viz.  $|\phi|$ . Hence, we have

$$\begin{aligned}
& \lim_{\varepsilon \downarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{r^0} (z_n \in B_\varepsilon(\psi)) \\
& \leq \lim_{\varepsilon \downarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{|r^0|} (z_n^r \in B_\varepsilon(|\psi|)) \\
& \leq \lim_{\varepsilon \downarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{|r^0|} (z_n^r \in B_\varepsilon(|\phi|)) \\
& = \lim_{\varepsilon \downarrow 0} \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{|r^0|} (z_n^r \in B_\varepsilon(|\phi|)) . \tag{7}
\end{aligned}$$

The last equality follows again from the LDP for the M/M/1 process.

We now claim that, in fact, both departure processes as well as the arrival process are nearly straight lines. This follows from [17, Lemma 7.25] as follows. Consider the augmented M/M/1 queue where we look at a 4-dimensional process, consisting of the 3 independent Poisson processes of the component process  $z^c$ , and the M/M/1 process whose arrival process is the Poisson  $\lambda$  process, and whose service process consists of (the superposition of) the other two independent processes. This M/M/1 process has exactly the desired distribution, and by [17, Lemma 7.25], each of the coordinate processes follows a straight line (note that in our case, the rates and hence  $\ell$  in [17, Eq. 7.25] do not depend on  $x$ , and  $r$  is linear so that  $y$  there is fixed, so that  $\theta(s)$  and  $\lambda_j(r(s))$  do not depend on  $s$ , and so  $wr^i$  there is a straight line). But if arrivals and departures follow a straight line, and the M/M/1 queue follows  $|\phi|$ , then by the definition of the JSQ, it follows  $\phi$ . Therefore, we have for all  $\varepsilon > 0$

$$\begin{aligned}
& \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{|r^0|} (z_n^r \in B_\varepsilon(|\phi|)) \\
& \leq \liminf_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{r^0} (z_n \in B_\varepsilon(\phi)) . \tag{8}
\end{aligned}$$

Take  $\varepsilon \downarrow 0$  in (8) and use (7) to conclude that  $\phi$  is more likely than  $\psi$ . ■

We now extend Lemma 6 to include the starting point 0.

**Lemma 7** *Let  $r^0 = (0, 0)$ ,  $r^T = (r_1, r_2)\beta$ . Define the constant speed path  $\phi(t) := (r_1, r_2)\beta t/T$ ,  $0 \leq t \leq T$ . Let  $\psi$  be another path on  $[0, T]$  from 0 to  $r^T$ , increasing strictly along the diagonal. Then  $\phi$  is more likely than  $\psi$ .*

*Proof.* The difficulty here is that if one of the queues is empty but the other is not, then  $z_n$  does not behave like  $z_n^r$ . Let  $t(\varepsilon)$  be the first time that  $\psi_i(t) > \varepsilon$  for all  $i$ . Then after  $t(\varepsilon)$ , if  $z_n$  is in  $B_\varepsilon(\psi)$  then none of the queues empties again. So let  $\phi_\varepsilon$  be the direct path from  $\psi(t(\varepsilon))$  to  $r^T$ . We can now repeat the arguments of Lemma 6, with two minor changes. First, here the starting point of  $z^r(t(\varepsilon))$  is  $|\psi(t(\varepsilon))|$  which is not necessarily equal

to  $|z_n(t(\varepsilon))|$ : however their distance for the paths of  $z_n$  that stay in  $B_\varepsilon(\psi)$  is at most  $\varepsilon$ . As  $\varepsilon \rightarrow 0$ , this does not change the proof. Finally,

$$\begin{aligned} \mathbb{P}_0(z_n \in B_\varepsilon(\psi)) &\leq \mathbb{P}_0(z_n \in B_\varepsilon(\psi), t \geq t(\varepsilon)) \\ &\leq \mathbb{P}_0(z_n^r \in B_\varepsilon(|\phi|), t \geq t(\varepsilon)) \end{aligned}$$

so that the argument of (7) applies once we note that

$$|\mathbb{P}_0(z_n^r \in B_\varepsilon(|\phi|)) - \mathbb{P}_0(z_n^r \in B_\varepsilon(|\phi|) \text{ for } t \geq t(\varepsilon))| \leq e^{-\eta(\varepsilon)n}$$

where  $\eta(\varepsilon) \rightarrow 0$  as  $\varepsilon \downarrow 0$ . With these changes the previous proof applies. ■

**Remark** The fact that paths with constant speed are most likely among all increasing paths along the diagonal also holds for decreasing paths. ■

Lemma 7 compares  $\phi$  only to other strictly increasing paths. We show below that for large  $T$  the most likely path stays near 0 and then follows  $\phi$ . We call a straight line path with constant speed a *direct path*. From the two preceding Lemmas we deduce

**Corollary 1** *For a direct path  $\phi$  increasing along the diagonal, the limit (2) exists. Let  $T$  be the duration of the path  $\phi$ , and  $\ell(\cdot)$  the local rate function of the (non empty)  $M(\lambda)/M(\mu)/1$  queue [17, Equation (7.17)]. Then*

$$I(\phi) = J(T; |\phi(T) - \phi(0)|) := T\ell(|\phi(T) - \phi(0)|/T), \quad (9)$$

where  $I$  is given in (1) and  $J$  is strictly convex in  $T$ . If  $\phi$  is decreasing than the result holds with a minus sign for  $|\phi(T) - \phi(0)|$  in (9).

*Proof.* That the limit exists follows directly from the proof of Lemma 6 (and of Lemma 7 in case the starting point is 0), particularly by inequalities (7) and (8). Also these inequalities say that the limit equals the Large Deviations cost rate going from  $|\phi(0)|$  to  $|\phi(T)|$  in  $T$  time units by the  $M(\lambda)/M(\mu)/1$  queue. The cost function is as stated in (9). Convexity follows by differentiation since  $\ell$  is strictly convex. ■

**Remark** A similar result (limit exists and exact expression) holds for any direct path  $\phi$  which does not cross the diagonal, for the following reason. Consider a (vector) process where the coordinates are statistically independent. Suppose an equation of the type (2) holds for each coordinate separately, with rate  $I_i(\phi_i)$  for the  $i$ th coordinate. Then it is easy to show that (2) also holds for the full process, and  $I(\phi) = \sum I_i(\phi_i)$ . More explicitly,

$$\begin{aligned} \{z_{1,n} \in B_{\varepsilon/2}(\phi_1) \text{ and } z_{2,n} \in B_{\varepsilon/2}(\phi_2)\} &\subset \{z_n \in B_\varepsilon(\phi)\} \\ &\subset \{z_{1,n} \in B_\varepsilon(\phi_1) \text{ and } z_{2,n} \in B_\varepsilon(\phi_2)\} \end{aligned}$$

so that the desired limit exists, and the rate functions add. Now a direct path that does not cross the diagonal is of one of the following types.

- Except possibly for the starting and/or ending point, the path lies entirely in one of the two interiors of the JSQ, i.e., either below or above the diagonal, and away from the boundaries. In these areas the system is equivalent to the local process, for which one coordinate is a Poisson process and the other an independent M/M/1 queue, so that we have an LDP (end points are easily dealt with in case they lie on a boundary, cf. proof of Lemma 7).
- The path goes along one of the boundaries. Here again the process is equivalent to the local process, comprising of a Poisson and an independent M/M/1 queue. Since the LDP holds for each of the independent coordinates, it also holds for the process, and the rate function is simply the sum of the two rates, one for each coordinate. ■

**Remark** We have not yet settled the optimal paths! Lemmas 6 and 7 say only that direct paths along the diagonal are the most likely increasing paths. The optimal path may decrease from the start until some time epoch and then increase. This was our first question earlier in this section. ■

First we consider paths starting at 0. When the time to overflow  $T$  is large enough, the optimal path remains at 0 and then moves straight with constant speed  $\mu_1 + \mu_2 - \lambda$  to the overflow point. This is due to the well-known fact that this is the most likely path to overflow in an M/M/1 queue with arrival rate  $\lambda$  and service rate  $\mu_1 + \mu_2$ —which is our reduced process. Otherwise, it leaves 0 immediately and moves with constant speed.

**Theorem 1** Let  $T_0^* := |r^T|(\mu_1 + \mu_2 - \lambda)^{-1}$ . The most likely way for the stable JSQ to reach the overflow point  $r^T = (r_1, r_2)\beta$  starting at 0, is:

- (i) If  $T \leq T_0^*$  then the optimal path is the direct path between 0 and  $r^T$ , i.e.,  $\phi(t) = (r_1, r_2)\beta t/T$ ,  $t \in [0, T]$ , with cost rate  $I(\phi) = J(T, |r^T|)$ .
- (ii) If  $T > T_0^*$  then the optimal path is to stay in 0 until  $T - T_0^*$ , and then proceed along the diagonal with speed  $\mu_1 + \mu_2 - \lambda$ , i.e.

$$\phi(t) = (r_1, r_2)(r_1 + r_2)^{-1}(\mu_1 + \mu_2 - \lambda)(t - (T - T_0^*)), \quad t \in [T - T_0^*, T],$$

and  $\phi(t) = 0$ ,  $t \in [0, T - T_0^*]$ , with cost rate

$$I(\phi) = J(T_0^*, |r^T|) = (r_1^T + r_2^T) \log((\mu_1 + \mu_2)\lambda^{-1}).$$

*Proof.* Let  $\psi$  a path containing a detour. That is, there are  $t_1 < t_2 < t_3$  with  $\psi(0) = 0$ ,  $\psi(T) = r^T$ ,  $\psi(t_1)$  and  $\psi(t_2)$  are on the diagonal and  $\psi(t)$  is not on the diagonal for  $t_1 < t < t_2$ . Since  $\min\{z_{1,n}(t)/r_1, z_{2,n}(t)/r_2\}$  can only increase on the diagonal, then if  $z_n$  is near  $\psi$  with positive probability, then necessarily  $\psi(t_2)$  is below and to the left of  $\psi(t_1)$ . Since  $\psi(T) = r^T$  and since increase only occurs on the diagonal, there must be a  $t_3$ , the smallest time after  $t_2$  so that  $\psi(t_3) = \psi(t_1)$ . Define  $\Delta = t_3 - t_1$  and consider

$$\bar{\psi}(t) = \begin{cases} 0 & 0 \leq t \leq \Delta \\ \psi(t - \Delta) & \Delta \leq t \leq t_3 \\ \psi(t) & t_3 \leq t. \end{cases} \quad (10)$$

We now use coupling to show that the probability to stay near  $\bar{\psi}$  is larger than that of  $\psi$ . To do this, consider the scaled component process  $z_n^c$ . Construct the coupled process  $\bar{z}_n^c$  by interchanging its segments as in (10). Since Poisson processes are memoryless, the distribution of the two processes is the same. By this coupling,

$$\begin{aligned} \mathbb{P}_0(z_n \in B_\varepsilon(\psi)) &\leq \mathbb{P}_0(\sup |z_n(t) - \psi(t)| < \varepsilon, \text{ for } t \in [0, t_1] \cup [t_3, T]) \\ &= \mathbb{P}_0(\sup |\bar{z}_n(t) - \bar{\psi}(t)| < \varepsilon, \text{ for } t \in [\Delta, T]). \end{aligned}$$

Since the JSQ system is stable,

$$\begin{aligned} \lim_{\varepsilon \downarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_0(\sup |\bar{z}_n(t) - \bar{\psi}(t)| < \varepsilon, \text{ for } t \in [\Delta, T]) \\ = \lim_{\varepsilon \downarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_0(\bar{z}_n \in B_\varepsilon(\bar{\psi})). \end{aligned}$$

Thus paths without detours are more likely. But then Corollary 1 applies, and the cost is that of the reduced M/M/1 queue, which gives (i) and (ii). ■

Theorem 1 and the analysis in the next sections fixes  $T$  a-priori. However, Definition 3 and the discussion below the definition relate events to sets of paths in a more general way. Let us illustrate how our results can be extended in the context of Theorem 1. Fix an overflow point  $r^F = \beta(r_1, r_2)$  and consider the event of reaching the overflow point from 0, without an a-priori restriction of the amount of time it takes (this is called a free time problem). That is, we are interested in both the probability as well as the optimal (most likely) path, including the most likely time the event takes.

Note that some care is required in formulating this question, since example,

$$\mathbb{P}_0(|z_n(T) - r^F| < \varepsilon \text{ for some } T) = 1. \quad (11)$$

To see this note that, since the system is stable, for each  $n$  the process will return to the empty state with probability one. Each time both queues are empty, the process “restarts,” independently of past behavior. We thus get repeated independent “experiments,” where in each (excursion from 0 back to 0) we have a very small, but non-zero probability of overflow. The Borell-Cantelli Lemma now implies the equality in (11).

To avoid this pitfall, let  $S$  be the following collection of paths:

$$S \triangleq \{\phi : \phi(0) = 0, \phi(t) \neq 0 \text{ for } t > 0 \text{ and } \phi(T_\phi) = r^F\}.$$

We claim that the large deviations limit of the probability of this event exists and that its rate is equal to  $J(T_0^*, |r^F|)$  mentioned in Theorem 1.

**Lemma 8** Fix  $T$  and let  $\phi_T^*$  be the most likely path of Theorem 1. Then

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_0(|z_n(T) - r^F| < \varepsilon) \\ = \lim_{\varepsilon \rightarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_0(z_n \in B_\varepsilon(\phi_T^*)) \\ = -I(\phi_T^*). \end{aligned}$$

*Proof.* This is in the spirit of [17, Lemma 2.8]. Consider the collection  $S'$  of all paths from 0 to  $r^F$  in  $T$  time units. The probability that the process stays close to any given path  $\phi \neq \phi_T^*$  is exponentially smaller than the corresponding probability of staying close to the optimal path  $\phi_T^*$ :

$$|\mathbb{P}_0(z_n \in B_\varepsilon(\phi_T^*)) - \mathbb{P}_0(z_n \in B_\varepsilon(\phi))| \leq e^{-n\gamma(\varepsilon)},$$

where  $\gamma(\varepsilon) = o(1)$  as  $\varepsilon \rightarrow 0$ . Using continuity of the rate function with respect to small shifts we obtain a similar result for paths  $\phi$  with  $|\phi(T) - r^F| < \varepsilon$ . Finally, since the process is exponentially tight, we can replace  $S'$  with a compact subset, and cover the compact subset with a finite collection of balls  $B_\varepsilon(\phi_i)$ . ■

The Lemma says that the event of overflow at a specific time happens most likely because the process follows its optimal path to overflow (upto the specified time). In the free time problem we optimize with respect to  $T$ . Let  $\phi^*$  be the optimal path, directly from 0 to  $r^F$  in  $T_0^*$  time units from Theorem 1. Because  $I(\phi^*) = J(T_0^*, |r^F|) = \inf_T I(\phi_T^*)$  we obtain

**Corollary 2**  $\lim_{\varepsilon \rightarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_0(z_n \in S) = -I(\phi^*).$

#### 4.1 General starting point

The methods leading to Theorem 1 can be used to derive the most likely path to overflow, starting at any point. Below we provide some results in this direction: for proofs we refer to [16]. We need the following fact.

**Lemma 9** *Let  $\phi$  be a concatenation of  $K$  direct paths. Then*

$$I(\phi) = -\lim_{\varepsilon \downarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} \log \mathbb{P}_{x^0}(z_n \in B_\varepsilon(\phi)) = \sum_{k=1}^K I(\phi_k). \text{ If } \phi_k \text{ does not cross the diagonal, then } I(\phi_k) \text{ has an explicit representation (e.g. (9)).}$$

The proof of this lemma uses standard continuity arguments, and in particular continuity of the rate function with respect to shifts, and homogeneity of the process away from the diagonal and the boundaries.

Optimal paths to overflow cannot be too complicated.

**Theorem 2** *The optimal path to overflow has the following properties. It is composed of a finite number of direct paths. The direction or speed do not change in the interiors (i.e. they only change either on the diagonal or when at least one coordinate equals 0). If it spends time below the diagonal, then it does not spend time above the diagonal and conversely. If  $\phi_i(t) = 0$  for some  $t > 0$ , then  $\phi_i(t) = 0$  in some interval containing  $t$ .*

The proofs of these statements use arguments as in Lemma 6. More can be obtained when the starting point is on the diagonal. Recall the definition of the reduced process, and the related cost function (9). The latter is the cost

rate for an M/M/1 queue, which is a strictly convex, unimodal function. Its minimum is easily determined:

$$T_d := \arg \min_{T>0} I_d(T) = |r^T - r^0|(\mu_1 + \mu_2 - \lambda)^{-1}.$$

**Theorem 3** *Suppose the starting point is on the diagonal.*

(i) *The direct path is optimal if  $T \leq T_d$ .*

(ii) *There is a  $T^*$  such that for  $T > T^*$ , the optimal path to overflow follows the most likely behavior to 0, stays there until  $T - T^*$  and then follows a direct path with speed  $\mu_1 + \mu_2 - \lambda$  to overflow. The cost rate equals*

$$(r_1^T + r_2^T) \log(\mu_1 + \mu_2) \lambda^{-1}.$$

It may be surprising at first glance that this is not the general case.

**Theorem 4** *Suppose the parameters are such that the most likely path moves down and away from the diagonal. Fix a non-zero starting point  $r^0$  on the diagonal. Then there is a  $T$  such that the most likely path to overflow in time  $T$  moves down until queue 2 empties, then moves left and then up to the diagonal, and finally through a direct path to overflow. Throughout the path, its first coordinate (queue 1) stays away from 0.*

For details and an explicit example see [16].

Note that once the most likely path is obtained for each  $T$ , the free time most likely path can be obtained as in Lemma 8 and Corollary 2.

The conclusions we can draw for the two-dimensional system are as follows. The most likely behavior can be computed explicitly, starting at any point. The most likely path to overflow can be found by solving a finite dimensional optimization problem: we need to compute the rate function over a finite collection of candidate paths, each composed of several direct path segments. For each such segment, the rate can be computed explicitly.

However, it is quite clear that the complexity of the problem grows quite fast as the dimension of the system increases. In the next section we comment on the extension.

## 5 Extension to higher dimensions

In this section we describe how the results for two dimensions can be extended to the more general model. As we shall see, the results for the most likely behavior extend naturally, and the complexity grows but is quite manageable. For the overflow problem the situation is more complex.

### 5.1 Most likely behavior in higher dimensions

The most likely behavior can be analyzed much in the same way as the two-dimensional system, and so we only comment on this.

(i) First note that if the JSQ system is unstable, namely  $\sum_{i=1}^K \mu_i < \lambda$  then, by the previous arguments, the most likely behavior is to approach the line  $r(t)$  and then follow that line, where the reduced system satisfies

$$z_{\infty}^r(t) = x^r + \left( \lambda - \sum_{i=1}^K \mu_i \right) \cdot t.$$

(ii) Suppose that some queues are empty, and some are not. Without loss of generality, assume that  $x_i = 0$ ,  $i = 1, \dots, i_x$  and  $x_i > 0$ ,  $i > i_x$ . Suppose moreover that the subsystem consisting of  $i_x$  queues is unstable, that is  $\sum_{i=1}^{i_x} \mu_i < \lambda$ . Then queues  $1, \dots, i_x$  are statistically independent of the rest of the queues until one of the other queues reaches a value of  $z_{i,n}(t)/r_i$  that is equal or lower than that of one of the first  $i_x$  queues. Therefore, for our local (short-time) analysis, we need not consider those queues: their behavior is simply to follow their drift  $-\mu_i e_i$ . We therefor ignore those queues, which amounts to setting  $K = i_x$ . Now by definition, the reduced system is unstable, and moreover, it is coupled to our system whenever all queues are non empty. But by definition, none of the ratios  $z_{i,n}(t)/r_i$  will be larger than 1 unless all are at least 1. Thus we see that the total size of the queues  $1, \dots, i_x$  increases at the rate  $\lambda - \sum_{i=1}^{i_x} \mu_i$  of the reduced system, and as analyzed before the increase is towards  $r$ .

(iii) If the subsystem is stable, than the same reasoning shows that these queues will remain empty. The queues we ignored receive no arrivals, and therefore grow smaller. As another queue empties, the stable system which now consists of one more queue is obviously stable (more departures with the same arrival rate), and so eventually all queues will empty.

(iv) Suppose now that there are no empty queues, but that a subset of the queues have nearly the same product  $x_i/r_i$ , and further, that this product is smaller than for other queues. Again, without loss of generality, assume  $x_i/r_i < x_j/r_j$  for all  $i \leq i_x < j$ . As before, we can consider the first group of queues, ignoring the rest. If the subsystem is unstable, that is  $\sum_{i=1}^{i_x} \mu_i < \lambda$  then, as before, the most likely path is for this set of queues to grow along  $r$ , until  $z_{i,n}(t)/r_i \geq z_{j,n}(t)/r_j$  for some  $j > i_x$ . This will eventually happen, as the other queues all decrease while the first queues all grow. At that point, if with the additional queue the system is still unstable, the present analysis holds, with new growth rates.

(v) Finally, if the subsystem is stable, it may decrease along  $r$ . To find out if this is the case, consider the queue (among the first  $i_x$  queues), say  $j$ , with the smallest ratio of  $\mu_i/r_i$ . This queue is our candidate to empty slower than the rest of the queues, thus staying “above the diagonal”  $r$ . To find out if this happens, we compute the rate at which the remaining queues empty. This is done by considering the reduced system with  $i_x - 1$  queues, then calculating the rate at which, say,  $z_{1,n}(t)/r_1$  decreases, and comparing this to  $\mu_j/r_j$ . If the latter is larger, then we know that queue  $j$  empties as fast as the remaining group, and the most likely behavior is for these  $i_x$  queues to decrease together along  $r$ . However, if  $\mu_j/r_j$  is smaller, than indeed it

will stay larger even though no arrivals join queue  $j$ . We then repeat the procedure, to check if another queue should be excluded.

Once this procedure is complete, The remaining group of queues empty along  $r$ , until they are all empty. This concludes the last case in our analysis.

We note that, unlike other systems, here the boundaries do not hamper our analysis. This is the case since, for our purposes, the coordinates decouple except at 0 and on  $r$ . In the first case, a one-dimensional analysis provides the picture. In the second case, the boundaries play no role, and the analysis near  $r$  can be completed.

### *5.2 Most likely path to overflow in higher dimensions*

We note that the methods of this paper apply to the higher dimensional system. These methods show that the rate function is the correct (asymptotic-logarithmic) value of the probabilities of overflow, even though the general large deviations principle may not hold. The main results, namely Theorem 1 as well as the results of subsection 4.1 hold as stated. Consequently, the infinite-dimensional optimization problem (finding the optimal path) is reduced to a finite dimensional problem: finding the parameters of the optimal path, which is constructed from a finite collection of straight lines with fixed speed. Moreover, the rate function for each such line possesses an explicit expression, since it is the rate for a collection of independent processes, where one process is an M/M/1 queue and the rest are Poisson processes. However, it is also clear that the complexity of the problem increases very fast with the dimension of the system, since the number of possibilities increases rapidly, and complicated paths, such as in Theorem 4 need to be considered. Fortunately, the behavior of this system is relatively simple since increase may only happen along a one-dimensional surface—the diagonal.

### **Acknowledgement**

The authors would like to thank an anonymous referee for valuable remarks and suggestions that improved the exposition.

### **References**

1. Atar, R. and P. Dupuis (1999) Large deviations and queueing networks: methods for rate identification. *Stoch. Proc. Appl.* 84: 255–296.
2. Bertsimas D., L. Paschalides and J.N. Tsitsiklis (1999) Large deviations analysis of the generalized processor sharing policy. *Queueing Systems* 32: 319–349.
3. Chang, C-S. (1995) Sample path large deviations andintree networks. *Queueing Systems* 20: 7–36.

4. Dembo, A. and O. Zeitouni (1996) *Large deviations techniques and applications*, Second edition. Springer Verlag.
5. Dupuis, P. and R.E. Ellis (1995) The large deviations principle for a general class of queueing systems I. *Trans. American Math. Soc.* 347: 2689–2751.
6. Dupuis, P. and R.E. Ellis (1995) *A Weak Convergence Approach to the Theory of Large Deviations*. Wiley.
7. Foley, R.D. and D. McDonald (2001) Join the shortest queue: stability and exact asymptotics. *Ann. Appl. Prob.* 11: 569–607.
8. Ganesh, A.J. and N. O’Connell (2002) A large deviation principle with queueing applications. *Stoch. and Stoch. Reports* 73: 25–35.
9. Ganesh, A.J., N. O’Connell and D. Wischik (2004) *Big Queues*. Springer.
10. Ignatiouk-Robert, I. (2005) Large deviations for processes with discontinuous statistics. *Ann. Prob.* to appear.
11. Mandjes, M. (1999) Rare event analysis of the state frequencies of a large number of markov chains. *Stoch. Models* 15: 577–592.
12. Mandjes, M. and A. Ridder (1999) Optimal trajectory in a queue fed by a large number of sources. *Queueing Systems* 31: 137–170.
13. Mandjes, M. and A. Weiss (2003) Sample path large deviations of a multiple time-scale queueing model. Preprint.
14. McDonald, D. (1996) Overloading parallel servers when arrivals join the shorest queue. pp. 169–196 In: Glasserman, P., K. Sigman and D. Yao (eds.) *Stochastic Networks: stability and rare events*. Springer Verlag.
15. Puhalskii, A.A. and W. Whitt (1998) Functional large deviation principles for waiting and departure processes. *Prob. Engin. Info. Sci.* 12: 479–507.
16. Ridder, A. and A. Shwartz (2005) Large deviations methods and the join-the-shortest-queue model. Report. Available (JSQreport) from <http://www.ee.technion.ac.il/~adam/PAPERS>.
17. Shwartz, A. and A. Weiss (1995) *Large deviations for performance analysis: queues, communication and computing*. Chapman Hall.
18. Shwartz, A. and A. Weiss (1999) Multiple time scales in Markovian ATM models I. Formal calculations. CC PUB 267, Elec. Engin. Technion.
19. Shwartz, A. and A. Weiss (2005) Large Deviations with diminishing rates. to appear, *Math. of Oper. Res.*.
20. Stolyar, A.L. (2003) Control of End-to-End Delay Tails in a Multiclass Network: LWDF Discipline Optimality. *Ann. Appl. Prob.* 13: 1151–1206.
21. Stolyar A.L. and K. Ramanan (2001) Largest Weighted Delay First Scheduling: Large Deviations and Optimality. *Ann. Appl. Prob.* 11: 1–48.
22. Turner, S. (2000) Large deviations for join the shorter queue. In: McDonald, D. and S. Turner (eds.) *Analysis of Communication Networks: Call Centres, Traffic and Performance*. The Fields Institute.
23. Wischik D.J. (2001) Sample path large deviations for queues with many inputs. *Ann. Appl. Prob.* 11: 379–404.
24. Weiss, A. (1986), A new technique for analyzing large traffic systems. *Adv. Appl. Prob.* 18: 506–532.