

THE LINEAR GEODESIC PROPERTY
IS NOT GENERALLY PRESERVED
BY A FIFO QUEUE¹

A. J. Ganesh² and Neil O'Connell³

Abstract

If a FIFO queue is fed by several input streams that jointly satisfy a sample path large deviation principle (LDP) with ‘linear geodesics’, then the cumulative departures (upto a large time) also satisfy the LDP with a rate function which depends in a relatively simple way on the rate function corresponding to the inputs: this was demonstrated in a recent paper by the second author. It suggests the possibility of an iterative scheme which would allow one to determine the large deviation behaviour of more complicated networks. To do this, however, one would require that the linear geodesic property be preserved: in this paper we demonstrate that in general it is not preserved. This is true even in the case of a single input stream.

¹Published in *Annals of Applied Probability*, 8(1), 98–111, 1998.

²Department of Statistics, Birkbeck College, University of London, Malet Street, London WC1E 7HX, U.K. E-mail: *A.Ganesh@stat.bbk.ac.uk* Research supported by a fellowship from BP and the Royal Society of Edinburgh.

³BRIMS, Hewlett Packard Laboratories, Filton Road, Stoke Gifford, Bristol BS12 6QZ, U.K. E-mail: *noc@hplb.hpl.hp.com*

1 Introduction and Preliminaries

There has been considerable recent interest in the large deviations behaviour of queueing systems. This started with the observation that for a single server queue, the tails of the queue length distribution can be characterised in terms of the large deviations behaviour of the arrivals and service processes. (This is actually a classical result, originally due to Cramér in the *iid* case; for more general statements in the context of queueing systems, see [2, 7, 8, 9].) Since then there have been many attempts to extend the theory to more complicated networks. A starting point in this quest is to consider the effect of interactions in a shared buffer, which is served according to a FIFO (first-in-first-out) policy. More precisely, if the arrival streams are assumed to jointly satisfy a large deviation principle (LDP), then what can be said about the joint large deviation behaviour of the corresponding departure streams? A partial answer to this question was presented in [7], where the notion of *decoupling of effective bandwidths* was introduced. There it is shown that there is a region over which the large deviation rate functions for the cumulative departures and arrivals agree, and bounds are given outside that region. Chang and Zajic [4] consider the case of a single arrival stream and stochastic service rate. In [10], a full description of the rate function for the cumulative departures is given in general, under the hypothesis that the arrival processes jointly satisfy a sample path LDP with ‘linear geodesics’ (roughly speaking, this means that the most likely path to an extreme value is a straight line⁴). This begs the question, which we will address in this paper: do the departures also satisfy this hypothesis? If so, then one could treat quite complicated networks by successive iteration of the single-buffer results in [10]. We know of one example where this is the case, namely if the inputs and service are independent Poisson processes and the queue is stable: then the outputs, in equilibrium, are also independent Poisson processes with the same rates as the corresponding inputs. However, we find that it is not generally the case, even when there is just one input stream.

The remainder of this section is devoted to giving some background and a formal description of the problem. Counterexamples to the above suggestion are given in the next section.

⁴It is usually assumed, or is a consequence of the stochastic model of the arrival process, that the sample paths of the arrival process satisfy an LDP with rate function of the form $I(\phi) = \int \Lambda^*(\dot{\phi}(t))dt$. If Λ^* is convex, then the sample path ϕ that minimizes $I(\phi)$ subject to the boundary conditions $\phi(0) = a$ and $\phi(1) = b$ is described by the straight line joining $(0, a)$ and $(1, b)$.

But, first, we discuss some special cases when the departure process does have linear geodesics. Suppose there is a single input stream. If the service process is deterministic, then the departure process has linear geodesics. So, a recursive analysis of networks of such queues is possible, as in [3]. Even if the service process is stochastic, we show that, conditional on the departure rate from a queue exceeding its mean, the departure process has linear geodesics. We are typically interested in the probability of queue lengths exceeding some large threshold, and in well-designed networks this requires departure rates exceeding their mean. Therefore, we have linear geodesics in the region of interest, and so the study of networks of queues using a recursive approach is again feasible. Such an approach has been taken in Bertsimas *et al.* [1], in the context of quite general arrival and service processes, and a single class of customers. We show in this paper that this approach can't be extended easily to networks with more than one traffic class. In fact, even if the service process is deterministic and there are only two traffic types, the departure process need not have linear geodesics. (This is true even if we condition on the aggregate departure rate exceeding its mean.)

We now give a formal description of the problem. Consider a discrete time queue with d arrival streams $\mathbf{X} = (X^1, \dots, X^d)$ sharing an infinite buffer according to an FIFO policy with stochastic service rate C . \mathbf{X}_k denotes the number of arrivals of each type in time slot k , while C_k denotes the maximum number of customers of any type that can be served in this time slot. We will begin by assuming that the queue is empty at time slot 0. Define

$$\mathbf{A}_n = \sum_{k=1}^n \mathbf{X}_k, \quad B_n = \sum_{k=1}^n C_k \quad (1)$$

Let $A_n = \sum_{j=1}^d \mathbf{A}_n^j$ denote the total number of arrivals, and D_n the total number of departures up to time n . Assuming that the queue is work-conserving, we have

$$D_n = \inf_{0 \leq k \leq n} (A_k - B_k) + B_n \quad (2)$$

The amount of work, $\mathbf{D}_n = (D_n^1, \dots, D_n^d)$, serviced from each input stream by time n , is defined as follows. Set,

$$T_n = \sup\{k \leq n : A_k \leq D_n\} \quad (3)$$

$$\mathbf{D}_n = \mathbf{A}_{T_n} + (D_n - A_{T_n})\mathbf{X}_{T_n+1}/X_{T_n+1}. \quad (4)$$

Recall that $\mathbf{X}_k = (X_k^1, \dots, X_k^d)$ denotes the amount of work arriving from the different streams in time slot k . T_n denotes the last time slot such that all arrivals up to it have been served by time n (some of the arrivals in time slot $T_n + 1$ may also have been served). In words, work is serviced in the order received and simultaneous arrivals from different sources are thoroughly mixed in the queue.

Define $\mathbf{S}_n(t) = (\mathbf{A}_{[nt]}/n, \mathbf{B}_{[nt]}/n)$, $\mathbf{R}_n(t) = \mathbf{D}_{[nt]}/n$. For each positive integer k , let \mathcal{L}^k denote the subspace of paths in $L_\infty([0, 1]^k)$ with non-decreasing components, and by $\mathcal{A}^k \subset \mathcal{L}^k$ the set of those paths with absolutely continuous components starting at zero. The following hypotheses are employed in [10].

(H1) For all $\gamma \in \mathbb{R}$, $\sup_k E[\exp \gamma(X_k + C_k)] < \infty$.

(H2) For each $\lambda \in \mathbb{R}^{d+1}$, the limit

$$\Lambda(\lambda) = \lim_{n \rightarrow \infty} \frac{1}{n} \log E[\exp(\lambda \cdot \mathbf{S}_n(1))] \quad (5)$$

exists as an extended real number and is finite in a neighborhood of the origin. The sequence \mathbf{S}_n satisfies the large deviation principle (LDP) in \mathcal{L}^{d+1} with good rate function I given by

$$I(\phi) = \begin{cases} \int_0^1 \Lambda^*(\dot{\phi}) ds, & \text{if } \phi \in \mathcal{A}^{d+1} \\ \infty, & \text{otherwise} \end{cases} \quad (6)$$

where Λ^* is the convex conjugate of Λ .

(H3) The arrival and service processes are asymptotically independent in the sense that

$$\Lambda^*(\mathbf{x}, c) = \Lambda_a^*(\mathbf{x}) + \Lambda_b^*(c) \quad (7)$$

We refer to the hypothesis (H2) as the ‘linear geodesic property’. It follows from (H2), the convexity of Λ^* and Jensen’s inequality, that the optimal path from point to point is a straight line. Such an LDP has been shown to hold quite generally by Dembo and Zajic [6]: roughly speaking, it holds provided the sequence is, in some sense, stationary and mixing. Under the above hypotheses, it was shown in [10] that the sequence \mathbf{R}_n satisfies the LDP in \mathcal{L}^d with good rate function given by

$$I_d(\psi) = \inf\{I(\phi) : \Delta(\phi) = \psi\} \quad (8)$$

where $\Delta : \mathcal{C}^{d+1} \rightarrow \mathcal{C}^d$ is defined by

$$\mathbf{A}(\phi) = (\phi^1, \dots, \phi^d) \quad (9)$$

$$D(\phi)(t) = \inf_{0 \leq \nu \leq 1} \left[A(\phi)(\nu t) - \phi^{d+1}(\nu t) \right] + \phi^{d+1}(t) \quad (10)$$

$$T(\phi)(t) = \inf\{r : A(\phi)(r) = D(\phi)(t)\} \quad (11)$$

$$\Delta(\phi) = \mathbf{A}(\phi) \circ T(\phi) \quad (12)$$

Here, (9) follows from (2), (10) from (3) and (11) from (4). $\mathbf{A}(\phi)$ denotes the (scaled) joint arrival process, while $A(\phi) = \phi^1 + \dots + \phi^d$ denotes the aggregate arrival process. For the scaled processes, $T(\phi)(t)$ denotes the last time, arrivals up to which depart by time t . Since the queue was assumed to be empty at time 0 and the service discipline is FIFO, the departures in all the streams up to time t , denoted $\Delta(\phi)(t)$, is precisely the arrivals in all the streams up to time $T(\phi)(t)$. $\Delta(\phi)$ describes the scaled departure process corresponding to the scaled arrival and service processes described by ϕ . By expressing the object of interest, the scaled departure process, as a continuous function, Δ , of the arrival and service processes, (12) sets the stage for applying the contraction principle. The contraction principle then yields the LDP in (8) for the *sample paths* of the departure process. Using the contraction principle once more, we obtain an LDP for the departure rate: \mathbf{D}_n/n satisfies the LDP in \mathbb{R}^d with good rate function

$$\Lambda_d^*(\mathbf{z}) = \inf\{I_d(\psi) : \psi(0) = 0, \psi(1) = \mathbf{z}\}, \quad (13)$$

for I_d as in (8). From this, it was derived in [10] that

$$\begin{aligned} \Lambda_d^*(\mathbf{z}) &= \inf\{\beta\Lambda_a^*(\mathbf{x}/\beta) + \sigma\Lambda_a^*\left(\frac{\mathbf{z} - \mathbf{x}}{\sigma}\right) + \beta\Lambda_b^*(c) + (1 - \beta)\Lambda_b^*\left(\frac{z - x}{1 - \beta}\right) : \\ &\quad \beta, \sigma \in [0, 1], c \in \mathbb{R}, \beta + \sigma \leq 1, x \leq \beta c\}. \end{aligned} \quad (14)$$

The last result has the interpretation that the most likely path of the arrival and service processes which results in the departure process having mean rate \mathbf{z} on the interval $[0, n]$ is as follows. The arrival process has rate \mathbf{x}/β on the interval $[0, \beta n]$ and rate $(\mathbf{z} - \mathbf{x})/\sigma$ on $[\beta n, (\beta + \sigma)n]$. The service rate during $[0, \beta n]$ is c , which is greater than the aggregate arrival rate during this period. So the queue is empty during $[0, \beta n]$. The queue is non-empty throughout $[\beta n, n]$, during which period the service rate, at $(z - x)/(1 - \beta)$, is no larger than the total arrival rate, which is $(z - x)/\sigma$. Therefore, the aggregate departure rate is equal to the aggregate arrival rate, x/β , during

the first phase, $[0, \beta n]$, when the queue is empty, and equal to the service rate, $(z - x)/(1 - \beta)$, during the second phase, $[\beta n, n]$, when the queue is never empty. The rigorous statement underlying this intuition, proved in [10], is the following:

$$\Lambda_d^*(\mathbf{z}) = I_d(\psi), \quad (15)$$

where $\psi(t)$, $0 \leq t \leq 1$ is specified by,

$$\psi(0) = 0, \quad \dot{\psi}(t) = \begin{cases} \frac{\mathbf{x}}{\beta}, & 0 < t < \beta, \\ \frac{z - \mathbf{x}}{1 - \beta}, & \beta < t < 1. \end{cases} \quad (16)$$

Here β , \mathbf{x} are those achieving the infimum in (14), and I_d is as defined in (8).

The result in (14) applies to a queue started empty. A similar but more involved expression was derived for a queue in equilibrium. Under the above hypotheses, we can derive an expression for the asymptotics of the queue length distribution. The problem of extending this derivation to an arbitrary queue in a feed-forward queueing network remains open. The arrival process into any queue in such a network is an aggregate of the departure processes from its predecessors (or splittings thereof) and possibly of an external arrival process. Therefore, the result above suggests that we approach this problem using the LDP for the departure process. This would work if the departure process also satisfied hypotheses (H1)-(H3). In the next section we give examples to show that there are situations where the departure process fails to satisfy (H2), both for the queue initially empty and for the system started in equilibrium.

2 Counterexamples

2.1 Single customer class

Consider a queue with a stochastic server and a single class of customers (so $d = 1$). Then, (14) simplifies to

$$\Lambda_d^*(z) = \inf\left\{\beta\Lambda_a^*(x/\beta) + \sigma\Lambda_a^*\left(\frac{z-x}{\sigma}\right) + \beta\Lambda_b^*(c) + (1-\beta)\Lambda_b^*\left(\frac{z-x}{1-\beta}\right) : \beta, \sigma \in [0, 1], c \in \mathbb{R}, \beta + \sigma \leq 1, x \leq \beta c\right\}. \quad (17)$$

where now x and z are scalars.

Let $EX = \Lambda_a'(0)$ denote the mean number of arrivals, and $EC = \Lambda_b'(0)$ the mean number of services in each time slot. The following properties of

Λ_a^* , Λ_b^* are well-known, see [5] for instance. Λ_a^* (respectively, Λ_b^*) is non-negative, and zero only at EX (respectively EC). Both Λ_a^* and Λ_b^* are convex, and finite on a non-empty interval, in the interior of which they are analytic. We assume that the queue is stable, namely $EX < EC$. Suppose that for some $\alpha \in [0, EX]$,

$$\Lambda_a^*(\alpha) = \Lambda_b^*(\alpha) \quad (18)$$

and also that these functions are finite in a neighborhood of α . Without loss of generality, we can take α to be the largest number in $[0, EX]$ for which (18) holds. Then, since $\Lambda_a^*(EX) = 0 < \Lambda_b^*(EX)$, we have $\Lambda_a^*(x) < \Lambda_b^*(x)$ for all $x \in (\alpha, EX]$, and consequently that $(\Lambda_a^*)'(\alpha) < (\Lambda_b^*)'(\alpha)$. It follows from this that

$$\exists \epsilon > 0 : \quad \Lambda_b^*(x) < \Lambda_a^*(x) < +\infty \quad \forall x \in [\alpha - \epsilon, \alpha], \quad (19)$$

and also that

$$\exists 0 < x_1 < \alpha < x_2 < EX : \quad \frac{x_1 + x_2}{2} = \alpha, \quad \frac{\Lambda_b^*(x_1) + \Lambda_a^*(x_2)}{2} < \Lambda_a^*(\alpha). \quad (20)$$

We shall show that in this case, the most likely departure path having mean rate α is not linear. Let $\phi(t) = (\phi^1(t), \phi^2(t))$, be defined on $[0, 1]$ by $\phi(0) = 0$ and

$$\dot{\phi}^1(t) = \begin{cases} x_2, & 0 < t < \frac{1}{2}, \\ EX, & \frac{1}{2} < t < 1. \end{cases} \quad \dot{\phi}^2(t) = \begin{cases} EC, & 0 < t < \frac{1}{2}, \\ x_1, & \frac{1}{2} < t < 1, \end{cases} \quad (21)$$

where x_1, x_2 are as in (20). Then, since $EX < EC$, we have from (9)-(12) that

$$\Delta(\phi)(t) = \begin{cases} x_2 t, & 0 \leq t \leq \frac{1}{2}, \\ \frac{1}{2} x_2 + (t - \frac{1}{2}) x_1, & \frac{1}{2} \leq t \leq 1. \end{cases} \quad (22)$$

and in particular that $\Delta(\phi)(1) = \alpha$. Therefore, by (8) and (13),

$$\begin{aligned} \Lambda_d^*(\alpha) &\leq I(\phi) \\ &= \int_0^1 \Lambda_a^*(\dot{\phi}^1(s)) ds + \int_0^1 \Lambda_b^*(\dot{\phi}^2(s)) ds \\ &= \int_0^{1/2} \Lambda_a^*(x_2) ds + \int_{1/2}^1 \Lambda_a^*(EX) ds + \int_0^{1/2} \Lambda_b^*(EC) ds + \int_{1/2}^1 \Lambda_b^*(x_1) ds \\ &= \frac{1}{2} [\Lambda_a^*(x_2) + \Lambda_b^*(x_1)] \\ &< \Lambda_a^*(\alpha) \end{aligned} \quad (23)$$

The first equality above follows from (6) and (7), the second from the definition of ϕ in (21), and the last from the fact that $\Lambda_a^*(EX) = \Lambda_b^*(EC) = 0$, see [5] for example. The last inequality above holds because of (20). Notice that the departure process $\Delta(\phi)$ in (22), corresponding to the arrival process ϕ^1 and service process ϕ^2 , is not linear but has different slopes x_2 and x_1 in two different periods of equal length.

Next, let $\psi(t)$ be linear on $[0, 1]$ with $\psi(0) = 0$ and $\psi(1) = \alpha$, so that $\dot{\psi}(t) = \alpha$ for all $t \in (0, 1)$. Consider any $\phi \in \mathcal{A}^2$ such that $\psi = \Delta(\phi)$. That is, (ϕ^1, ϕ^2) is any pair of arrival and service processes (excluding those whose rate function is $+\infty$), corresponding to which ψ is the departure process. Then, by (9)-(12),

$$\psi(t) = \inf_{0 \leq s \leq t} [\phi^1(s) - \phi^2(s)] + \phi^2(t), \quad (24)$$

from which it is clear that $\psi(t) \leq \phi^1(t)$. In particular, $\phi^1(1) \geq \alpha$. If $\phi^1(1) = \alpha$, then, by (6),

$$\begin{aligned} I(\phi) &\geq \int_0^1 \Lambda_a^*(\dot{\phi}^1(s)) ds \\ &\geq \Lambda_a^* \left(\int_0^1 \dot{\phi}^1(s) ds \right) \\ &= \Lambda_a^*(\alpha). \end{aligned} \quad (25)$$

The first inequality is due to the non-negativity of Λ_a^* , Λ_b^* , the second holds because of Jensen's inequality and the convexity of Λ_a^* , while the equality is because $\phi^1(1)$ was assumed to be α . If $\phi^1(1) > \alpha$, define

$$\tau = \sup\{t \in [0, 1] : \phi^1(t) \leq \alpha t\} \quad (26)$$

and note that $\tau < 1$. Hence, by continuity of ϕ , $\phi^1(\tau) = \alpha\tau$.

Lemma 1 *Suppose that $\psi(t) = \alpha t$ for all $t \in [0, 1]$, where ψ is defined by (24), and that $\phi \in \mathcal{A}^2$. Then, with τ given by (26),*

$$\phi^2(t) - \phi^2(\tau) = \alpha(t - \tau) \quad \forall t \in [\tau, 1].$$

Proof : As noted above, $\phi^1(\tau) = \alpha\tau = \psi(\tau)$, the latter equality holding by hypothesis regarding ψ . From this, we see that the infimum in (24) corresponding to $t = \tau$ is achieved at $s = \tau$. Consequently, (24) implies that

$$\psi(t) = \inf_{\tau \leq s \leq t} [\phi^1(s) - \phi^2(s)] + \phi^2(t) \quad \forall t \geq \tau. \quad (27)$$

If the infimum above is achieved at τ for all $t \in [\tau, 1]$, then, for all t in this interval, $\phi^2(t) - \phi^2(\tau) = \psi(t) - \psi(\tau) = \alpha(t - \tau)$, and so the lemma is established. Otherwise, because ϕ is absolutely continuous, one of the following must hold:

$$\exists \epsilon > 0 : \dot{\phi}^1(s) - \dot{\phi}^2(s) < 0 \quad \forall s \in (\tau, \tau + \epsilon), \quad (28)$$

or

$$T \triangleq \inf\{s > \tau : \phi^1(s) - \phi^2(s) < \phi^1(\tau) - \phi^2(\tau)\} \in (\tau, 1). \quad (29)$$

In the former case, the infimum in (27) corresponding to $t = \tau + \epsilon$ is achieved at $s = \tau + \epsilon$, and so

$$\psi(\tau + \epsilon) - \psi(\tau) = \phi^1(\tau + \epsilon) - \phi^1(\tau) > \alpha\epsilon, \quad (30)$$

where the inequality follows from the definition of τ in (26). In the latter case, we see from the continuity of ϕ that the infimum in (27) corresponding to $t = T$ is achieved at $s = T$. So $\psi(T) = \phi^1(T)$ and

$$\psi(T) - \psi(\tau) = \phi^1(T) - \phi^1(\tau) > \alpha(T - \tau), \quad (31)$$

where the inequality follows from (26). Now, both (30) and (31) contradict the hypothesis that $\psi(t) = \alpha t$ for all $t \in [0, 1]$. Therefore, neither (28) nor (29) can hold, implying that the infimum in (27) must be achieved at τ for all $t \in [\tau, 1]$. This completes the proof of the lemma. □

From the above lemma and (26), we obtain using (6) that

$$\begin{aligned} I(\phi) &\geq \int_0^\tau \Lambda_a^*(\dot{\phi}^1(s)) ds + \int_\tau^1 \Lambda_b^*(\dot{\phi}^2(s)) ds \\ &\geq \tau \Lambda_a^*(\alpha) + (1 - \tau) \Lambda_b^*(\alpha) \\ &= \Lambda_a^*(\alpha). \end{aligned} \quad (32)$$

The first inequality is due to the non-negativity of Λ_a^* and Λ_b^* , and the second is due to their convexity and Jensen's inequality (note that $\phi^1(\tau) = \alpha\tau$, while $\phi^2(1) - \phi^2(\tau) = \alpha(1 - \tau)$ by Lemma 1). The equality follows from the definition of α in (18).

Let ψ be given by $\psi(t) = \alpha t$, $t \in [0, 1]$. Since either (25) or (32) applies to any $\phi \in \mathcal{A}^2$ for which $\Delta\phi = \psi$, observe from (8) that $I_d(\psi) \geq \Lambda_a^*(\alpha)$. Therefore, by (23), ψ does not achieve the infimum in (13) corresponding

to $z = \alpha$. In other words, the departure process with constant rate α is not the most likely to achieve an average departure rate α ; this is achieved by a process with a nonlinear path. This implies that I_d cannot be expressed in the form

$$I_d(\xi) = \int_0^1 \Lambda^*(\dot{\xi}) ds$$

for any convex function Λ^* and so the departure process does not satisfy hypothesis (H2).

The above conclusion applies to a queue started empty. We now consider a queue in stationarity. Let Q_0 denote the queue length at time 0. It is shown in [10] that the scaled queue lengths Q_0/n satisfy an LDP in \mathbb{R} with rate function L , which is explicitly computed. For our purposes, it is enough to note that $L(0) = 0$ and that $L(q) \geq 0$ for all $q > 0$. Suppose the scaled initial queue length is q , and that the scaled process of arrivals and services is described by $\phi = (\phi^1, \phi^2)$. Then, the scaled departure process up to time t is given by

$$D(q, \phi)(t) = \phi^2(t) \wedge \inf_{0 \leq \nu \leq 1} [q + \phi^1(\nu t) - \phi^2(\nu t) + \phi^2(t)], \quad (33)$$

where $x \wedge y$ denotes $\min\{x, y\}$. Notice that since $\phi(0) = \mathbf{0}$, we recover (10) for the departure process from an empty queue by substituting $q = 0$. We shall show that for any $q > 0$ and $\phi \in \mathcal{L}^2$ such that $D(q, \phi(t)) = \alpha t$ for all $t \in [0, 1]$, we have $L(q) + I(\phi) \geq \Lambda^*(\alpha)$. This will enable us to conclude that the departure process does not have linear geodesics, even in equilibrium.

Let ψ be linear on $[0, 1]$ with $\psi(t) = \alpha t$ for all $t \in [0, 1]$. Fix $q > 0$ and let $\phi \in \mathcal{C}^2$ be such that $D(q, \phi) = \psi$. Then, by (33), either $\phi^2(t) = \alpha t$ for all $t \in [0, 1]$, or $q + \phi^1(s) - \phi^2(s) < 0$ for some $s \in [0, 1]$. In the former case, we have by (6), (7) and the non-negativity of the Λ^* that

$$I(\phi) \geq \int_0^1 \Lambda_b^*(\dot{\phi}^2(s)) ds = \Lambda_b^*(\alpha). \quad (34)$$

In the latter case, we have by the continuity of ϕ that

$$\tau \triangleq \inf\{s \in [0, 1] : q + \phi^1(s) - \phi^2(s) < 0\} \in [0, 1).$$

It follows that $D(q, \phi(t)) = \phi^2(t)$ for all $t \in [0, \tau]$, whereas, for $t \in [\tau, 1]$,

$$D(q, \phi)(t) - D(q, \phi)(\tau) = \phi^2(t) - \phi^2(\tau) + \inf_{\tau \leq s \leq t} [q + \phi^1(s) - \phi^2(s)],$$

because $q + \phi^1(s) - \phi^2(s)$ takes its minimum value on $[0, \tau]$ at τ , and this value is zero. Hence, we can rewrite the above as

$$\begin{aligned} & D(q, \phi)(t) - D(q, \phi)(\tau) \\ &= \inf_{\tau \leq s \leq t} \left[\left(\phi^1(s) - \phi^1(\tau) \right) - \left(\phi^2(s) - \phi^2(\tau) \right) \right] + \phi^2(t) - \phi^2(\tau). \end{aligned} \quad (35)$$

Define $\tilde{\phi}(t) = \phi(t) - \phi(\tau)$, $t \in [\tau, 1]$. Then, we have from above that

$$\psi(t) = \begin{cases} \phi^2(t), & \text{if } t \in [0, \tau], \\ \inf_{\tau \leq s \leq t} \left[\tilde{\phi}^1(s) - \tilde{\phi}^2(s) \right] + \tilde{\phi}^2(t) + \phi^2(\tau), & \text{if } t \in (\tau, 1]. \end{cases} \quad (36)$$

Comparing this with (24), we see that the departure process on $[\tau, 1]$ is identical to that from an empty queue with arrival and service processes given by $\tilde{\phi}$. This is not surprising because the queue does, in fact, become empty at time τ by definition of τ . Since $\tilde{\phi}$, restricted to $[\tau, 1]$, is merely a shifted version of ϕ on this interval, $I(\phi) = I(\tilde{\phi})$ for ϕ restricted to this interval. Therefore,

$$I(\phi) \geq \int_0^\tau \Lambda_b^*(\dot{\phi}^2(s)) ds + I(\tilde{\phi}).$$

Now, since $\psi(t) = \alpha t$, $\dot{\phi}^2(s) = \alpha$ for all $s \in [0, \tau]$. Also, by the same derivation as leads to (25) and (32), we have $I(\tilde{\phi}) \geq (1 - \tau)\Lambda_a^*(\alpha)$. Finally, since $\Lambda_a^*(\alpha) = \Lambda_b^*(\alpha)$ by definition of α , and $L(q) \geq 0$ for all q , we get

$$L(q) + I(\phi) \geq \Lambda_a^*(\alpha).$$

This holds for all initial queue lengths $q \geq 0$, and arrival and service processes ϕ , that result in a linear departure process $\psi(t) = \alpha t$. Note that (23) continues to hold for departures in equilibrium because it was derived for departures from an empty queue, and we have $L(0) = 0$, see [10]. Therefore,

$$\Lambda_d^*(\alpha) < \inf\{L(q) + I(\phi) : D(q, \phi)(t) = \alpha t \ \forall t \in [0, 1]\}, \quad (37)$$

which implies that, conditional on a mean departure rate of α , the most likely path is not linear. Thus, even in equilibrium, the departure process does not necessarily have linear geodesics.

We end this subsection with some comments about the scope and implications of the above results. A careful look at the proof shows that the result relied on α being less than EX and on the rate functions of the arrival

and service processes intersecting at α . If the service process is deterministic, the latter cannot happen, and in this case it can be shown that the departure process has linear geodesics. This makes it possible to analyze networks of deterministic server queues, as in Chang [3]. Likewise, if we consider only $\alpha > EX$, then too it can be shown that the departure process conditioned on having mean rate α is linear. Since we are typically interested in the problem of queue lengths exceeding some large threshold, and since in well-designed networks this requires departure rates exceeding their mean, we are usually only interested in the rate function of departures for $\alpha > EX$. Since we have linear geodesics in this region, the study of networks of queues using a recursive approach is again feasible. Such an approach has been taken in Bertsimas *et al.*, [1]. We shall next show that neither of these features comes to our rescue when dealing with multi-class queues. In this case, the joint departure process can have non-linear geodesics even if the server is deterministic, and even if we consider departures whose aggregate rate exceeds their mean.

2.2 Two customer classes

Consider a queue multiplexing two customer classes, and served deterministically at rate c . Suppose that customers from the first class arrive deterministically at rate a , while those of the second have a stochastic arrival process satisfying hypotheses (H1)-(H3) with the rate function Λ_2^* . We assume that the mean aggregate arrival rate is strictly less than the service rate, c . Note that the two arrival streams are trivially independent, as are the arrival and service processes. We denote the large deviations rate function of the first arrival process by Λ_1^* and that of the service process by Λ_b^* . So

$$\Lambda_1^*(x) = \begin{cases} 0, & \text{if } x = a, \\ +\infty, & \text{else.} \end{cases} \quad \Lambda_b^*(x) = \begin{cases} 0, & \text{if } x = c, \\ +\infty, & \text{else.} \end{cases} \quad (38)$$

For some $\epsilon > 0$ and $b < c + \epsilon - a$, let $\mathbf{z} = (a - \epsilon, b)$ and consider the departure process conditioned to have mean rate \mathbf{z} . We shall show that this departure process does not have linear geodesics.

Let $\psi \in \mathcal{A}^2$ be linear with $\psi(0) = 0$ and $\psi(1) = \mathbf{z}$, so $\dot{\psi} = \mathbf{z}$. We show that there is no $\phi \in \mathcal{A}^3$ with $I(\phi) < +\infty$ such that $\Delta(\phi) = \psi$. In other words, there is no process of arrivals and services whose rate function is finite, corresponding to which ψ is the departure process. Suppose otherwise. Let $\phi \in \mathcal{A}^3$ have $I(\phi) < +\infty$, so that

$$\phi^1(t) = at, \quad \phi^3(t) = ct, \quad (39)$$

and suppose that

$$\Delta(\phi)(t) = \psi(t) = ((a - \epsilon)t, bt), \quad (40)$$

where $\epsilon > 0$ and $a + b - \epsilon < c$. Observe from (9)-(12) that

$$\Delta(\phi)(t) = \left(\phi^1(T(\phi)(t)), \phi^2(T(\phi)(t)) \right). \quad (41)$$

Therefore, by (39) and (40), $T(\phi)(t) = (a - \epsilon)t/a$, and so $\phi^2((a - \epsilon)t/a) = bt$. In addition, by (11),

$$D(\phi)(t) = \phi^1(T(\phi)(t)) + \phi^2(T(\phi)(t)) = (a + b - \epsilon)t. \quad (42)$$

But, by (10),

$$D(\phi)(t) = \inf_{0 \leq s \leq t} [\phi^1(s) + \phi^2(s) - \phi^3(s)] + \phi^3(t),$$

and so, by (39) and the fact, noted above, that $\phi^2(s) = abs/(a - \epsilon)$, we get

$$\begin{aligned} D(\phi)(t) &= \inf_{0 \leq s \leq t} \left[as + \frac{abs}{a - \epsilon} - cs \right] + ct \\ &= \begin{cases} ct, & \text{if } \frac{a}{a - \epsilon}(a + b - \epsilon) \geq c, \\ \frac{a}{a - \epsilon}(a + b - \epsilon)t, & \text{else.} \end{cases} \end{aligned}$$

Because of our hypothesis that $a + b - \epsilon < c$, we have $D(\phi)(t) > (a + b - \epsilon)t$ in either case above, contradicting (42). We have thus shown that, if $I(\phi) < +\infty$, then $\Delta(\phi) = \psi$ is impossible for $\psi(t) = \mathbf{z}t$ with $\mathbf{z} = (a - \epsilon, b)$. Therefore, by (8), $I_d(\psi) = +\infty$.

We now show that $\Lambda_d^*(\mathbf{z}) < +\infty$ for \mathbf{z} as above. Since $\epsilon > 0$ was arbitrary, we assume without loss of generality that $a - 2\epsilon > 0$ and define

$$x_1 = a + 2(b - \epsilon) - c, \quad x_2 = \frac{a}{a - 2\epsilon}(c - a + 2\epsilon). \quad (43)$$

Since the only requirement we imposed above was that $a + b - \epsilon < c$, it is clear that b and ϵ can be chosen so that $x_1 \geq 0$. Also, $x_2 > 0$ since it was assumed that c is larger than a . Let $\phi \in \mathcal{A}^3$ be defined by

$$\phi(0) = 0, \quad \dot{\phi}(t) = \begin{cases} (a, x_1, c), & 0 < t < 1/2, \\ (a, x_2, c), & 1/2 < t < 1. \end{cases} \quad (44)$$

Since x_1 and x_2 are non-negative, ϕ has non-decreasing components as required by the definition of \mathcal{A}^3 . Note that

$$a + x_1 = 2(a + b - \epsilon) - c < c$$

by the hypothesis that $a + b - \epsilon < c$, whereas

$$a + x_2 = \frac{ac}{a - 2\epsilon} > c.$$

In other words, the aggregate arrival rate $a + x_1$ is less than the service rate c during $[0, 1/2]$ whereas, at $a + x_2$, it is greater than c during $[1/2, 1]$. Therefore, the joint departure process $\Delta(\phi)$ is given by

$$\Delta(\phi)(0) = 0, \quad \frac{d}{dt}\Delta(\phi)(t) = \begin{cases} (a, x_1), & 0 < t < 1/2, \\ \left(\frac{a}{a+x_2}c, \frac{x_2}{a+x_2}c\right), & 1/2 < t < 1. \end{cases} \quad (45)$$

This is intuitively clear from the description of the queue, but can also be formally established using (9)-(12). Hence, we have from (43) that

$$\begin{aligned} \Delta(\phi)^1(1) &= \frac{1}{2} \left(a + \frac{ac}{a+x_2} \right) = a - \epsilon, \\ \Delta(\phi)^2(1) &= \frac{1}{2} \left(x_1 + \frac{x_2c}{a+x_2} \right) = b. \end{aligned}$$

Therefore, by definition of \mathbf{z} , $\Delta\phi(1) = \mathbf{z}$. Furthermore, by (44), (6) and (38), we have

$$\Lambda^*(\phi) = \frac{1}{2} (\Lambda_2^*(x_1) + \Lambda_2^*(x_2)) \quad (46)$$

for x_1, x_2 as in (43). Therefore, $\Lambda^*(\phi)$ is finite if $\Lambda_2^*(x_1)$ and $\Lambda_2^*(x_2)$ are, as is true if, for instance, the second arrival process is Poisson. It now follows from (8) and (13) that $\Lambda_d^*(\mathbf{z}) < +\infty$. But we showed earlier that $I_d(\psi) = +\infty$ for ψ given by $\psi(t) = \mathbf{z}t$. Therefore, the departure process with linear path does not achieve the infimum in (13), implying that the departure process does not satisfy a large deviations principle with action functional that is the integral of a convex rate function. In other words, it is not true that

$$I_d(\xi) = \int_0^1 \Lambda^*(\xi(s)) ds$$

for any convex function Λ^* . Consequently, the joint departure process does not satisfy hypothesis (H2), and so a recursive approach to estimating asymptotics of the queue lengths in a network does not appear feasible.

We now consider the same queueing system in stationarity, rather than started empty. It is shown in [10] that in stationarity, the scaled queue lengths \mathbf{Q}_0/n satisfy an LDP in \mathbb{R}^2 with a rate function L that can be computed explicitly. Here $\mathbf{Q}_0 = (Q_0^1, Q_0^2)$ denotes the number of customers

of each of the two types in the queue at time zero. It suffices for our purposes to note that $L(\mathbf{q}) \geq 0$ for all $\mathbf{q} \geq \mathbf{0}$, with equality if $\mathbf{q} = \mathbf{0}$.

Consider the system starting at time zero with scaled queue length $\mathbf{Q}_0/n = \mathbf{q}$. Suppose the arrival and service processes are given by $\phi = (\phi^1, \phi^2, \phi^3)$, and that $\psi = (\psi^1, \psi^2)$ is the corresponding departure process. Let a and c be defined as above to be the deterministic rate of the first arrival process and the service process respectively. Let $\epsilon \in (0, a)$ and $b > 0$ be such that $a + b - \epsilon < c$. We shall show that if $\psi(t) = \mathbf{z}t$ for all $t \in [0, 1]$, where $\mathbf{z} = (a - \epsilon, b)$, then $L(\mathbf{q}) + I(\phi) = \infty$.

Analogous to (33), the scaled process of aggregate departures up to time t is given by

$$D(\mathbf{q}, \phi)(t) = \phi^3(t) \wedge \inf_{0 \leq \nu \leq 1} [q + A(\phi)(\nu t) - \phi^3(\nu t) + \phi^3(t)], \quad (47)$$

where $q \triangleq q^1 + q^2$ is the total number in queue at time zero, and $A(\phi) \triangleq \phi^1 + \phi^2$ is the aggregate arrival process. Note that setting $q = 0$ above recovers (10) since $\phi(0) = \mathbf{0}$. Now, if ψ is to be the departure process, then we must have $D(\mathbf{q}, \phi) = \psi^1 + \psi^2$. Since $\psi(t) = \mathbf{z}t$, with $\mathbf{z} = (a - \epsilon, b)$, the above implies that $D(\mathbf{q}, \phi) = (a + b - \epsilon)t$ for all $t \in [0, 1]$. Recall that if $I(\phi)$ is to be finite, then we must have $\phi^1(t) = at$ and $\phi^3(t) = ct$ for all $t \in [0, 1]$, since the first arrival process and the service process are deterministic with rates a, c respectively. Therefore, for all such ϕ , (47) implies that

$$(a + b - \epsilon)t = ct \wedge \inf_{0 \leq s \leq t} [q + A(\phi)(s) - cs] + ct.$$

Since $a + b - \epsilon < c$, the above implies that $\inf_{0 \leq s \leq t} [q + A(\phi)(s) - cs]$ is strictly negative for all $t > 0$. Now $A(\phi)(s) = \phi^1(s) + \phi^2(s)$, $\phi(0) = \mathbf{0}$ and ϕ is continuous if $I(\phi)$ is finite. Therefore, it follows from the above that $\mathbf{q} = \mathbf{0}$, i.e., the queue must start empty. Then, by the argument above for departures from an empty queue, there is no process ϕ of arrivals and services such that the departure process is ψ and $I(\phi)$ is finite. We have also shown that this conclusion does not change if we allow any positive initial queue size, \mathbf{q} . This completes the proof that $I_d(\psi) = \infty$ even in equilibrium, where I_d denotes the rate function of the departure process.

We argued above that $\Lambda_d^*(\mathbf{z}) < \infty$ for departures from an empty queue. Since $L(0) = 0$, see [10], this argument applies to departures in equilibrium as well. Thus, the most likely path leading to a mean departure rate \mathbf{z} is not linear. This also implies that the rate function $I_d(\psi)$ for equilibrium departures cannot be of the form $\int_0^1 \Lambda^*(\dot{\psi})(s) ds$, for any convex function Λ^* .

Therefore, the departure process does not satisfy Hypothesis H2, needed to apply the results of [10] inductively to feed-forward multi-class queueing networks.

3 Conclusion

We considered the problem of characterising the large deviations behaviour of the departure process from a FIFO queue multiplexing several traffic streams. Such a characterisation could, in principle, be used iteratively to determine the large deviations behaviour of all processes of interest in networks of queues, and thereby to obtain the tail of the queue length and waiting time distributions at each queue in the network. The starting point of our analysis was the general description, in [10], of the rate function for the cumulative departures as the solution of a variational problem. It was shown in [10] that if, in addition, the arrivals satisfy a ‘linear geodesics’ condition, then the variational problem reduces to a finite-dimensional optimization problem. Such a simplification is essential if an iterative approach to analysing networks of queues is to be practical. This naturally leads to the question of whether the departures also satisfy the ‘linear geodesics’ assumption. We showed in this paper that this is not generally the case, even when there is just one input stream. Nevertheless, in the case of a single input stream, the departures do satisfy the linear geodesics requirement in the regime leading to large queue sizes. So an iterative approach to obtaining the tail of the queue length is possible, see [1]. However, such is not the case for multiple traffic streams, even when the service rate is deterministic.

References

- [1] D. Bertsimas, I. Paschalidis and J. Tsitsiklis (1994). On the large deviations behaviour of acyclic networks of G/G/1 queues. *Preprint*.
- [2] Cheng-Shang Chang (1994). Stability, queue length and delay of deterministic and stochastic queueing networks. *IEEE Trans. on Automatic Control* 39:913–931.
- [3] C.S. Chang (1995). Sample path large deviations and intree networks. *Queueing Systems*, 20(1-2): 7-36.

- [4] C.-S. Chang and T. Zajic (1995). Effective bandwidths of departure processes from queues with time varying capacities. *IEEE INFOCOM Proceedings*.
- [5] A. Dembo and O. Zeitouni (1993). *Large Deviations Techniques and Applications*, Jones and Bartlett.
- [6] A. Dembo and T. Zajic (1995). Large deviations: from empirical mean and measure to partial sums process. *Stoch. Proc. Appl.* 57:191–224.
- [7] G. de Veciana, C. Courcoubetis and J. Walrand (1994). Decoupling bandwidths for networks: a decomposition approach to resource management. *IEEE INFOCOM Proceedings*.
- [8] N.G. Duffield and Neil O’Connell (1995). Large deviations and overflow probabilities for the general single server queue, with applications. *Proc. Camb. Phil. Soc.* 118(1).
- [9] Peter W. Glynn and Ward Whitt (1995). Logarithmic asymptotics for steady-state tail probabilities in a single-server queue. *J. Appl. Prob.*, to appear.
- [10] Neil O’Connell (1997). Large deviations for departures from a shared buffer. *J. Appl. Prob.*, to appear.