

# Many-Sources Large Deviations for Max-Weight Scheduling

Vijay G. Subramanian, Tara Javidi, and Somsak Kittipiyakul

## Abstract

In this paper, a many-sources large deviations principle (LDP) for the transient workload of a multi-queue single-server system is established where the service rates are chosen from a compact, convex and coordinate-convex rate region and where the service discipline is the max-weight policy. Under the assumption that the arrival processes satisfy a many-sources LDP, this is accomplished by employing Garcia's extended contraction principle that is applicable to quasi-continuous mappings.

For the simplex rate-region, an LDP for the stationary workload is also established under the additional requirements that the scheduling policy be work-conserving and that the arrival processes satisfy certain mixing conditions.

The LDP results can be used to calculate asymptotic buffer overflow probabilities accounting for the multiplexing gain, when the arrival process is an average of *i.i.d.* processes. The rate function for the stationary workload is expressed in term of the rate functions of the finite-horizon workloads when the arrival processes have *i.i.d.* increments.

## Index Terms

max-weight policy, many-sources LDP, quasi-continuity, Garcia's extended contraction principle.

Vijay Subramanian is with the Hamilton Institute, National University of Ireland, Maynooth, Co. Kildare, Ireland. He acknowledges SFI for the support of this research through grant 07/IN.1/I901.

Tara Javidi is with the Department of Electrical and Computer Engineering, University of California, San Diego, CA, USA.

Somsak Kittipiyakul is with the School of Information, Computer and Communication Technology, Sirindhorn International Institute of Technology, Thammasat University, Thailand. This work was carried out as a part of his Ph.D. dissertation while at the Department of Electrical and Computer Engineering, University of California, San Diego, CA, USA.

## I. INTRODUCTION

The drive to achieve maximum efficiency in wireless data networks and high-speed switches has led to many advances in the design of good scheduling policies. One such family of good scheduling policies is an online policy<sup>1</sup> called the maximum weight (max-weight) scheduling policy. For the typical multi-class queue where only one queue can be served at a time, the max-weight policy serves one of the queues that has the largest value for the product of the workload and the service rate. We are interested in applying the max-weight policy to wireless networks where the scheduler is able to change the operating parameters at different levels of the traditional networking stack. Thus, the server has access to a richer choice of service options when compared to a traditional multi-class queue setting: each service option is a point within a compact, convex and coordinate-convex rate region. In this setting, the max-weight policy naturally generalizes to finding an operating point within the rate region that has the maximum projection along the workload vector<sup>2</sup>. Note that in this setting the traditional single-server multi-class queue has a rate-region given by a simplex.

In the present article, with  $K$  independent queues we seek to derive the probability of buffer overflow, when the server scheduling follows a max-weight policy. More specifically, for a given finite value  $B > 0$ , we consider the two buffer overflow quantities. First, we consider  $P(\mathbf{W}_{0,T} \geq B\mathbf{1}_K)$  where  $\mathbf{W}_{0,T} \in \mathbb{R}_+^K$  is the transient workload (to be formally defined later) at time 0 with “zero” initial workload at time  $-T$  and  $\mathbf{1}_K \in \mathbb{R}_+^K$  is the vector of all 1s. The second quantity we study is the stationary overflow probability for the limiting workload vector as  $T \rightarrow \infty$ , i.e.,  $P(\mathcal{W} \geq B\mathbf{1}_K)$ . Since these probabilities are, in general, very hard to compute exactly, we consider logarithmic asymptotics to the probabilities of interest using the theory of large deviations. In particular, this paper, under a “many-sources” scaling regime, establishes logarithmic asymptotics for 1) the transient workload for a compact, convex, and coordinate-convex rate region; and 2) the stationary workload for a simplex rate region using a work conserving scheduler.

<sup>1</sup>Online policies are those that can only use the past history of the arrivals, workloads and service decisions to decide on the scheduling choice; for example, these policies are not even aware of average arrival rates.

<sup>2</sup>In many ways cross-layer optimization has resulted in some firm strides towards a union of information theory and communication networks of the sort that was sought in [11].

In the classical<sup>3</sup> many-sources asymptotic, one considers a sequence of queueing systems indexed by the number of (independent) sources multiplexed (or averaged) over a particular queue, i.e., the arrival process to each queue is the average of  $L$  processes. The analysis focuses on the asymptotic behavior of the systems when  $L \rightarrow \infty$ . In our work, we consider a generalization of many sources asymptotic in which the input to the queueing system  $L$  exhibits a sample path large deviations property (LDP) similar to that of the average of  $L$  independent arrival streams (See Assumption 1). Given a sample path large deviations principle (see Definition 3) for the arrival processes, we derive a large deviations principle for the workload under the max-weight scheduling polic. In particular, we first show that the finite-horizon workload is a quasi-continuous map of the arrival process, for both the regular version of the max-weight policy and for a work-conserving version of it. Then the first contribution of the paper is that the finite-horizon workloads satisfy an LDP. This is obtained using a recent extension of the contraction principle by J. Garcia [41]. Restricting our attention to the simplex rate region (corresponding to the traditional multi-class single server queue), we again use Garcia’s extended contraction principle (along with a mixing condition assumption on the arrival process) to establish an LDP for the stationary workload. We should emphasize here that in contrast to related “many-sources” LDP results on FCFS and Priority policies that can be shown to be continuous, our LDP is established for an inherently discontinuous map that results from the max-weight scheduling policy. The LDP results (Theorems 1, 2 and 3) directly imply that the probability of buffer overflow has an exponential tail whose decay rate is dictated by a good rate function determined by the statistics of the arrival process. This rate function can be expressed as a solution to a finite-dimensional optimization problem which has the same flavor of a deterministic optimal control problem. The final contribution of our work is to provide a simplified form for the corresponding rate functions, when the arrival process has *i.i.d.* increments.

The outline of the paper is as follows. In Section II, we briefly motivate and contextualize our work in the larger body of literature on LDP analysis of queues as well as cross layer scheduling. The problem formulation is given in Section III. Section IV provides background and preliminary results. The main results of the paper, which are the LDPs of the workloads,

<sup>3</sup>The appellation “classical” is taken from [53] where a general scaling framework is presented that encompasses in a single setting all the different scalings used in the “many-sources” scaling regime.

are given in Section V and proved in Section VI. Section VII gives simplified expressions of the rate functions. We conclude in Section VIII with a discussion of future work.

We close this section with a summary of various notation used in this work. We use bold letters to discriminate vectors from scalar quantities as well as their components. We denote the set of natural and non-negative real numbers by  $\mathbb{N}$  and  $\mathbb{R}^+$ , respectively. We take  $\otimes$  to represent the Kronecker product. For  $0 \leq m_1 \leq m_2$  integers and a vector sequence  $(\mathbf{A}_t, t \in \mathbb{N})$  where  $\mathbf{A}_t \in \mathbb{R}_+^K$  for  $K \in \mathbb{N}$ , we define  $\mathbf{A}(m_1, m_2] := \sum_{t=m_1+1}^{m_2} \mathbf{A}_t$  as the cumulative arrivals from  $m_1 + 1$  until timeslot  $m_2$  where addition applies coordinate-wise. For vector-valued sequence  $\mathbf{A}$  we write  $\mathbf{A}|_{(m_1, m_2]}$  to denote the finite subsequence  $\{\mathbf{A}_{-m_2}, \dots, \mathbf{A}_{-m_1-1}\}$ . For a vector  $\mathbf{x} \in \mathbb{R}_+^K$  and set  $B \subset \mathbb{R}_+^K$ ,  $\text{Proj}_B(\mathbf{x})$  denotes the projection of vector  $\mathbf{x}$  on the set  $B$ ,  $\text{int}(B) := \{\boldsymbol{\lambda} \in B : \exists \boldsymbol{\lambda}' \in B \text{ s.t. } \boldsymbol{\lambda} < \boldsymbol{\lambda}'\}$  denotes the set of points strictly inside  $B$ , and  $[\mathbf{x}]^+ := \max\{\mathbf{0}, \mathbf{x}\}$  where the function applies coordinate-wise. Lastly, for any given function  $F : \mathcal{X} \mapsto \mathcal{Y}$  on metric spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , and  $x \in \mathcal{X}$ , we use the notation

$${}^x F := \{y \in \mathcal{Y} : (\exists x_n \rightarrow x) \text{ such that } F(x_n) \rightarrow y\},$$

to denote the set of all cluster points (in  $\mathcal{Y}$ ) of the images of sequences in  $\mathcal{X}$  converging to a point  $x \in \mathcal{X}$ . Note that  $(\cdot)F$ , in general, is a correspondence (also called a set-valued function) from  $\mathcal{X}$  to  $\mathcal{P}(\mathcal{Y})$  (the power set of  $\mathcal{Y}$ ). However,  $(\cdot)F$  is single-valued at  $x$ , i.e.,  ${}^x F = \{F(x)\}$ , if and only if  $F$  is continuous at  $x$ .

## II. RELATED WORK

In recent years, cross-layer scheduling has become a major focus of research in queueing and information theory due to its potential applications in communication networks. For brevity we do not list many important works, and instead mention only those closely related to this work.

In our LDP analysis, we follow the lead of many recent papers on the analysis of scheduling algorithms [9], [10], [12]–[16], [18] by considering logarithmic asymptotics to the probabilities of certain rare events. Maximum weight scheduling policy falls under class of the generalized  $c\mu$ -rule policies and is known to be stabilizing under very mild conditions [1]–[5]. A refined analysis of this policy shows that it minimizes the workload in the heavy traffic regime [6]–[8] over a large class of stationary online policies. This optimality of the max-weight policies also carries over to Large Deviations based tail asymptotes: the work-conserving version of these

policies is known to minimize the exponent of the tail asymptote of the stationary workload over a large class of stationary, online and work-conserving policies [10].

The present paper is closely related to [14], [15], where the buffer overflow probability for the workload processes of a single-server multi-queue queueing system under max-weight policies and general compact, convex and coordinate-convex capacity regions was established. While [14], [15] addresses the “large-buffer” scaling regime, this paper establishes similar logarithmic asymptotics results under the “many-sources” scaling regime (see [16]–[26], [30]–[39]). As the body of work on the “many-sources” scaling regime has grown, results have been established for many different scheduling policies for a single-server queue and also for networks of queues, namely, FCFS [20]–[22], [24], [26]–[28], [30]–[33], [38], [39], priority queueing [18], [24], [26], [36], [38], GPS [37], and SRPT and similar policies [16], [17], [34]. Our LDP analysis of max-weight scheduling is strongly motivated and complemented by these papers.

Finally, we close this section with a discussion of our motivation to consider the “many-sources” scaling studied in this paper. The interest in the “many-sources” asymptotic is known to be best motivated by 1) a recent and practical interest in applications when there are large number of flows to each user or node. This asymptote usually gives a more refined approximation to the probabilistic quantities of interest by incorporating the impact of the multiplexing gain [16]–[39] obtained by averaging many traffic sources together. However, our interest in the “many-sources” asymptotic has also been fueled by our earlier work [40] on 2) a cross-layer optimization of the PHY layer parameters, e.g., duration of the finite code blocks or cooperative cluster size when the fading channel is operated at high signal-to-noise-ratio (SNR). The high SNR regime is a very natural setting for the many-sources scaling since the capacity of the channel typically scales to infinity as  $\log(\text{SNR})$ , and therefore it is natural to scale the arrival rate of the flows with the same parameter, which is best accomplished by multiplexing more sources; in other words by setting  $L \propto \log(\text{SNR})$ . The present work on the many sources large deviation analysis of max-weight provides a first step in extending the above cross-layer optimization to a multi-user setting, an important topic for future research.

### III. PROBLEM FORMULATION

We consider a discrete-time queueing system with  $K \in \mathbb{N}$  independent queues and one server. We are interested in the statistical properties of the unfinished workload in queue  $k$  at time  $t$

under a max-weight server allocation policy. Let  $W_t^k \in \mathbb{R}_+$  be the unfinished workload (queue length) of queue  $k$  at the beginning of time  $-t$  and  $R_t^k$  be the amount of service allocated to queue  $k$  during time  $(-t, -t+1]$ . Let  $\mathbf{W}_t := (W_t^k, k \in \mathcal{K})$  be the corresponding workload vector and  $\mathbf{R}_t := (R_t^k, k \in \mathcal{K})$  be the rate vector. For every queue  $k \in \mathcal{K} := \{1, \dots, K\}$  we assume that work (in bits) arrives into the queue given by a sequence  $(A_t^k, t \in \mathbb{N})$  where  $A_t^k \in \mathbb{R}_+$  is the work brought in at time  $-t$ . For  $t \in \mathbb{N}$ , the dynamics of the workloads of queue  $k \in \mathcal{K}$  is

$$W_{t-1}^k = [W_t^k - R_t^k]^+ + A_t^k. \quad (1)$$

Note that we assume that the arrivals  $A_t$  happen any time in  $(-t, -t+1)$  but cannot be served in that timeslot  $-t$ .

The set of server's operating points is restricted to a compact and convex set  $\mathcal{R} \subset \mathbb{R}_+^K$  known as the *capacity* or *rate region* of the server, i.e.,  $\mathbf{R}_t \in \mathcal{R}$ , for all time  $t$ . We make the simplifying assumption that bits are infinitely divisible so that the rate allocations can be assumed to be real numbers. Furthermore, we also assume that  $\mathcal{R}$  is coordinate-convex, i.e, if  $\boldsymbol{\mu}_1 \in \mathcal{R}$ , then every  $\boldsymbol{\mu}_2 \in \mathbb{R}_+^K$  such that  $\boldsymbol{\mu}_2 \leq \boldsymbol{\mu}_1$  is also in  $\mathcal{R}$  where the inequalities apply along each coordinate. We are interested in the max-weight scheduler, and its closely related work-conserving version. At the beginning of timeslot  $-t$ , the rate vector  $\mathbf{R}_t \in \mathcal{R}$  is selected by a max-weight scheduler in response to the current workload  $\mathbf{W}_t$ . Specifically, under max-weight scheduler and in response to the current workload  $\mathbf{W}_t$ , the rate vector  $\mathbf{R}_t^*$  is chosen such that

$$\mathbf{R}_t^* \in \arg \max_{\mathbf{R} \in \mathcal{R}} \langle \mathbf{R}, \mathbf{W}_t \rangle. \quad (2)$$

As later established by Lemma 1, it is possible to construct a quasi-continuous (see Defn 1 in Section IV) function  $H$  such that  $\mathbf{R}_t^* = H(\mathbf{W}_t)$ . We call this construction the (max-weight) scheduling function. We also define a non-idling modification of max-weight scheduler for which the rate of service  $\mathbf{R}_t^{**}$  is such that it splits the service when the unfinished workload in each queue  $k$  is less than  $C^k = \max\{R^k : \mathbf{R} \in \mathcal{R}\}$ . Lemma 1 also shows that it is possible to construct a quasi-continuous function  $H^{\text{wc}}$  such that  $\mathbf{R}_t^{**} = H^{\text{wc}}(\mathbf{W}_t)$  and

$$H^{\text{wc}}(\mathbf{W}_t) = \begin{cases} \text{Proj}_{\mathcal{R}}(\mathbf{W}_t) & \text{if } \mathbf{W}_t \in \prod_{k=1}^K [0, C_k); \\ H(\mathbf{W}_t) & \text{otherwise.} \end{cases} \quad (3)$$

When necessary we will distinguish the workload vectors that result from the work-conserving max-weight by labeling them as  $\mathbf{W}^{\text{wc}}$ .

As mentioned earlier, in this paper we are interested in the probability distributions for the *finite-horizon* and *infinite-horizon* workloads. The finite-horizon workload, denoted by  $\mathbf{W}_{0,T}$ , is the workload at time 0, assuming the initial condition at time  $-T$  is  $\mathbf{W}_T = \mathbf{0}$ . The index  $T$  in  $\mathbf{W}_{0,T}$  reminds us of this initial condition.<sup>4</sup> The infinite-horizon workload,  $\mathcal{W}$ , is defined as  $\mathcal{W} := \lim_{T \rightarrow \infty} \mathbf{W}_{0,T}$ . We assume that the limit exists but may be infinite. Note that our results for the infinite-horizon workload are obtained in the restricted setting of a work conserving max-weight scheduler operating on a simplex rate region. For this work-conserving max-weight scheduler, it is known that  $\mathcal{W}^{\text{wc}}$  is the stationary workload when the system is stable.

We will use functions  $G_T$  and  $G_T^{\text{wc}}$  to relate the arrival process  $\mathbf{A}|_{(0,T]}$  with the unfinished work under max-weight scheduling,  $\mathbf{W}_{0,T}$  and under work conserving max-weight  $\mathbf{W}_{0,T}^{\text{wc}}$ , i.e.  $\mathbf{W}_{0,T} = G_T(\mathbf{A}|_{(0,T]})$  and  $\mathbf{W}_{0,T}^{\text{wc}} = G_T^{\text{wc}}(\mathbf{A}|_{(0,T]})$ . Similarly, we also define function  $G^{\text{wc}}$  to describe the workload under work-conserving max-weight when the arrival sequence is given, i.e.  $\mathcal{W}^{\text{wc}} = G^{\text{wc}}(\mathbf{A})$ ; we do not indicate the rate-region here as it is implicitly understood to be the simplex rate-region.

For each user  $k \in \mathcal{K}$  and system indexed by  $L \in \mathbb{N}$ , we will assume a stationary arrival process of work brought into the system given by a sequence  $A^{k,L} := (A_t^{k,L}, t \in \mathbb{N})$  where  $A_t^{k,L} \in \mathbb{R}_+$  is the work (in bits) brought in at time  $-t$  into the queue of user  $k$ . The arrivals to different queues/users are assumed to be mutually independent. Also let  $\mathbf{A}^L := (A^{k,L}, k \in \mathcal{K})$  be the sequence of arrival vectors. In our large deviation analysis, we characterize the asymptotic probability distributions for the finite-horizon and infinite-horizon workloads,  $G_T(\mathbf{A}^L|_{(0,T]})$  and  $G_T^{\text{wc}}(\mathbf{A}^L|_{(0,T]})$ , as  $L \rightarrow \infty$ .

We close this section by noting that in its typical avatar a max-weight scheduler is also accompanied by a set of non-zero weights  $\beta \in \mathbb{R}_+^K$  to weigh the workload vectors while determining the max-weight service vector. Using the observations in [14], [15] it can be seen that there is no loss of generality in assuming that the weight vector is  $\mathbf{1}_K$ .

<sup>4</sup>Note that the result remains valid even when the initial condition is within  $\mathcal{R}$  with the work-conserving scheduler. With  $\mathbf{W}_T \in \mathcal{R}$ , we always have the workload at time  $-T+1$  be  $\mathbf{W}_{T-1} = [\mathbf{W}_T - H^{\text{wc}}(\mathbf{W}_T)]^+ + \mathbf{A}_T = \mathbf{A}_T$  from the non-idling condition that we imposed on the server allocation mechanism as  $\text{Proj}_{\mathcal{R}}(\mathbf{W}_T) = \mathbf{W}_T$ .

#### IV. BACKGROUND AND PRELIMINARIES

In this section, we provide a brief review of fundamental definitions, concepts, and relevant results in large deviations theory that are essential for understanding our paper. Except for Lemma 1, which establishes the quasi-continuity of the scheduling functions  $H$  and  $H^{\text{wc}}$ , the material in this section can be found in [24], [26], [41], [51].

##### A. Quasi-continuity and Almost Compactness

In this section, we recall two important analytic properties for functions on metric spaces: quasi-continuity and almost compactness. These properties allow for an extension of contraction principle to which we later appeal. First, let us provide the definition of quasi-continuity on metric spaces:

*Definition 1:* [41, Theorem 3.2] Let  $\mathcal{X}, \mathcal{Y}$  be complete metric spaces. Function  $F : \mathcal{X} \mapsto \mathcal{Y}$  is *quasi-continuous* at  $x \in \mathcal{X}$  if and only if for each  $x \in \mathcal{X}$ , there is a sequence  $\{x_n\}$  such that  $x_n \rightarrow x, F(x_n) \rightarrow F(x)$ , and such that for all  $n$ ,  $F$  is continuous at  $x_n$ .

*Remark 1:* Obviously, every continuous function is quasi-continuous. A step function  $F : \mathbb{R} \mapsto \mathbb{R}$ , where  $F(x) = 0$  for  $x < 0$ ,  $F(x) = 1$  for  $x \geq 0$ , is quasi-continuous. However, if  $F(0) = 1/2$ , then  $F$  is not quasi-continuous.

*Remark 2:* An important property that we will use later is that if  $F$  is a continuous function and  $G$  is a quasi-continuous function, then  $F \circ G$  is quasi-continuous. However,  $G \circ F$  is not necessarily quasi-continuous [41].

As stated before, max-weight and its work conserving variation allow for quasi-continuous function selections.

*Lemma 1:* There exist quasi-continuous functions  $H$  and  $H^{\text{wc}}$  such that:

$$H(\mathbf{W}_t) \in \arg \max_{\mathbf{R} \in \mathcal{R}} \langle \mathbf{R}, \mathbf{W}_t \rangle,$$

and

$$H^{\text{wc}}(\mathbf{W}_t) = \begin{cases} \text{Proj}_{\mathcal{R}}(\mathbf{W}_t) & \mathbf{W}_t \in \Pi_{k=1}^K [0, C_k) \\ H(\mathbf{W}_t) & \text{otherwise} \end{cases}.$$

*Proof:* See Appendix A. The proof relies on the structural properties of the scheduling maps. ■

Having provided the definition of quasi-continuity, we are ready to define almost compactness for a function:

*Definition 2:* [41, Lemma 6.1] If  $\mathcal{X}, \mathcal{Y}$  are complete metric spaces, a function  $F : \mathcal{X} \mapsto \mathcal{Y}$  is *almost compact* at  $x \in \mathcal{X}$  if for every sequence  $x_n$  converging to  $x$ , there is a subsequence along which  $F$  converges to a point  $y \in \mathcal{Y}$ . We say that  $F$  is almost compact if it is almost compact at every  $x \in \mathcal{X}$ .

### B. Topology for Sample Paths

Since a large deviations principle is defined using topological entities and since we will deal with continuity and convergence of the workload mappings, we need to precisely specify the topology for the space of the arrival sample paths. We use the scaled-uniform norm/weighted supremum norm topology <sup>5</sup> used in [24] for our analysis.

Let  $\mathcal{D} \subset \{x : \mathbb{N} \mapsto \mathbb{R}_+\}$  denote the space of non-negative sequences such that  $\sup_{t \in \mathbb{N}} \left| \frac{x(0,t)}{t} \right| < +\infty$ , and let  $\mathcal{D}^K$  be the  $K$  cartesian product of  $\mathcal{D}$ . Let  $\|\cdot\|_u$  be the scaled uniform norm on  $\mathcal{D}$ , i.e.,  $\|x\|_u := \sup_{t \in \mathbb{N}} \left| \frac{x(0,t)}{t} \right|$  for all  $x \in \mathcal{D}$ , while for all  $a = (a^k, k \in \mathcal{K}) \in \mathcal{D}^K$ , where  $a^k \in \mathcal{D}$ , the scaled uniform norm of  $a$  is  $\|a\|_u := \max_{k \in \mathcal{K}} \|a^k\|_u$ . The space  $\mathcal{D}$  is metrizable via the scaled uniform norm  $\|\cdot\|_u$ , i.e., for all  $x, y \in \mathcal{D}$ , the distance between them is  $\|x - y\|_u = \sup_{t \in \mathbb{N}} \left| \frac{x(0,t) - y(0,t)}{t} \right|$ . Define a subspace  $\mathcal{D}_\mu$  of  $\mathcal{D}$  which contains all the arrival paths whose average arrival rate is equal to the expected rate  $\mu$ , i.e.,  $\mathcal{D}_\mu := \left\{ x \in \mathcal{D} : \lim_{t \rightarrow \infty} \frac{x(0,t)}{t} = \mu \right\}$ . Also for a vector  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$  define  $\mathcal{D}_\mu^K$  to be the product space of  $\mathcal{D}_{\mu^k}$  for all  $k \in \mathcal{K}$ . We equip  $\mathcal{D}_\mu$  and  $\mathcal{D}_\mu^K$  with the appropriate subspace and product topologies [42]. For finite dimensional metric spaces like  $\mathbb{R}_+^n$ ,  $n \in \mathbb{N}$ , we use the square uniform topology (same as the product topology) with the square metric  $\rho$  [42], where  $\rho(\mathbf{x}, \mathbf{y}) := \max_{i \in \{1, \dots, n\}} |x^i - y^i|$ . From [24], [26] it is also clear that the scaled uniform norm topology is stronger than the point-wise convergence topology, hence the projection and shift operators are continuous under the scaled uniform norm topology.

<sup>5</sup>In the theory of weak convergence of probability measures this topology was first proposed in [43], [45] for continuous functions vanishing at infinity with a continuous index set. In the same context it was then generalized to *cad-lag* functions with a continuous index set in [44]. The most general setting of this topology for discrete-time processes can be found in [46], [47], and the corresponding setting for continuous-time processes can be found in [48]. The usage of this topology in the context of Large Deviations can be found in [49] and [50]. Finally, the central theme in [24], [26], [50] is to demonstrate how this is a natural topology to use in the queueing context.

### C. Large Deviations Principle

The following definition of a large deviations principle is taken from [24]. For a complete reference to the theory, definitions, and tools, see [51].

*Definition 3 (Large deviations principle):* A sequence of random variables  $X^L$  in a Hausdorff space  $\mathcal{X}$  with  $\sigma$ -algebra  $\mathcal{B}$  is said to satisfy a large deviations principle (LDP)<sup>6</sup> with good rate function  $I$  if, for any  $B \in \mathcal{B}$ ,

$$-\inf_{x \in B^\circ} I(x) \leq \liminf_{L \rightarrow \infty} \frac{1}{L} \log P(X^L \in B) \leq \limsup_{L \rightarrow \infty} \frac{1}{L} \log P(X^L \in B) \leq -\inf_{x \in \bar{B}} I(x), \quad (4)$$

where  $B^\circ$  and  $\bar{B}$  are the interior and the closure of  $B$ , respectively, and if the rate function  $I : \mathcal{X} \mapsto \mathbb{R}_+ \cup \{\infty\}$  has compact level sets, where the level sets are defined as  $\{x : I(x) \leq \alpha\}$ , for  $\alpha \in \mathbb{R}$ .

If  $X^L$  is a mapping from  $\mathbb{N}$  to  $\mathbb{R}$  describing the sample path of a random sequence, the LDP is referred to as a *sample path* LDP.

### D. Garcia's Extended Contraction Principle

The contraction principle (see [51, p. 126]) says that if we have an LDP for a sequence of random variables, we can effortlessly obtain LDP's for a whole other class of random sequences that are obtained via continuous transformations. However, due to the inherent discontinuity in the max-weight scheduling function, this (regular) contraction principle fails to provide sufficient structure in the setting of our interest. Instead, we will utilize the following powerful extension of the contraction principle for quasi-continuous transformations on metric spaces, given by Garcia [41]. Garcia's extended contraction principle [41, Theorem 1.1] then says the following:

*Fact 1:* Assume that  $\Omega \xrightarrow{X^L} \mathcal{X} \xrightarrow{F} \mathcal{Y}$ ,  $\mathcal{X}, \mathcal{Y}$  are metric spaces, and  $\{X^L\}$  satisfies an LDP in  $\mathcal{X}$  with rate function  $I^\sharp$ . If at every  $x$  with  $I^\sharp(x) < \infty$ , the following hold:

- 1)  $F$  is almost compact; and
- 2) for all  $y \in {}^x\mathbb{F}$ , there exists a sequence  $\{x_n\}$  converging to  $x$  such that  $F(x_n) \rightarrow y$ ,  $F$  is continuous at  $x_n$ , and  $I^\sharp(x_n) \rightarrow I^\sharp(x)$ ,

then  $\{F(X^L)\}$  satisfies an LDP with rate function given by

$$I(y) = \inf \{I^\sharp(x) : y \in {}^x\mathbb{F}\}. \quad (5)$$

<sup>6</sup>Often  $\mathcal{B}$  is taken to be the Borel  $\sigma$ -algebra, and a rate function is, by definition, non-negative and lower semicontinuous.

*Remark 3:* Whenever  $I^\sharp(\cdot)$  is continuous, then the second condition is reduced to confirming that  $F$  is a quasi-continuous function.

## V. ASSUMPTIONS AND OVERVIEW OF THE RESULTS

Garcia's extended contraction principle together with Lemma 1 suggests the following road map to obtaining an LDP for the finite and infinite horizon workload processes under max-weight scheduling. The large deviation property for the sequences of finite- and infinite-horizon workloads would follow as a direct consequence of the sample-path LDP of the arrival process, as soon as one establishes the quasi-continuity and almost compactness of the mappings  $G_t$ ,  $G_t^{\text{wc}}$ , and  $G^{\text{wc}}$  along with some continuity properties of the rate function obtained from the sample-path LDP assumption on the arrival processes. As a result, our first task is to restrict our attention to arrival streams that satisfy the sample-path LDP as stated by Assumptions 1-3 in Section V-A. There we also discuss a family of arrival processes which satisfies these assumptions.

### A. Sample Path LDP of Arrival Processes

The following sample path LDP for the sequence of arrival processes  $\mathbf{A}^L$  is the starting point of our analysis.

*Assumption 1 (Many-sources sample path LDP):* The sequence  $\{\mathbf{A}^L\}$  satisfies a sample path LDP in  $\mathcal{D}_\mu^K$  equipped with the scaled uniform topology with rate function  $I^\sharp$ , where the rate function  $I^\sharp$  is given as

$$I^\sharp(\mathbf{a}) := \sup_{t \in \mathbb{N}} I_t^\sharp(\mathbf{a}|_{(0,t]}) = \lim_{t \rightarrow \infty} I_t^\sharp(\mathbf{a}|_{(0,t]}) \quad (6)$$

for  $\mathbf{a} \in \mathcal{D}_\mu^K$ , where for every  $t \in \mathbb{N}$  we also assume that  $\{\mathbf{A}^L|_{(0,t]}\}$  satisfies an LDP with rate function  $I_t^\sharp(\cdot)$  (in the product topology).

*Remark 4:* The most general conditions for Assumption 1 to be satisfied are given in [24, Thm. 3] (also stated in [26, Thm. 7.1, pg. 156]). There it is also shown how several standard stationary processes used for traffic modeling, such as *i.i.d.* increment processes, Markov-modulated, a general class of Gaussian, and fractional Brownian processes (for long-range dependent or heavy-tailed traffic), satisfy Assumption 1. The conditions of [24, Thm. 3] also imply that the sequence  $\{\mathbf{A}^L\}$  also satisfies an LDP on  $\mathcal{D}^K$  equipped with the scaled uniform topology, with rate function  $I^\sharp$  where  $I^\sharp(a) = \infty$  for  $a \in \mathcal{D}^K / \mathcal{D}_\mu^K$  [24], [26]. Finally, it is shown in [26, Lemma 7.8] that

under the conditions of [24, Thm. 3], for all  $t \in \mathbb{N}$  we have  $I_t^\sharp(\boldsymbol{\mu} \otimes \mathbf{1}_t) = 0$  and also that  $I_t^\sharp(\cdot)$  is convex.

In this paper, we further assume the following continuity conditions on the rate functions:

*Assumption 2:* We assume that  $I_t^\sharp(\cdot)$  is continuous in the product topology on  $\mathfrak{R}_+^{K \times t}$ .

*Assumption 3:* For every point  $x$  in the effective domain of  $I^\sharp$ , i.e.,  $\{\mathbf{a} \in \mathcal{D}_\mu^K : I^\sharp(\mathbf{a}) < +\infty\}$ , we assume that for every sequence  $\{\mathbf{a}^n\}$  converging to  $\mathbf{a}$  in  $\mathcal{D}_\mu^K$  such that there exists  $t \in \mathbb{N}$  so that  $\mathbf{a}_s^n = \mathbf{a}_s$  for all  $s > t$  (for all  $n$ ), we have  $I^\sharp(\mathbf{a}^n) \rightarrow I^\sharp(\mathbf{a})$ .

*Remark 5:* Assumption 3 is a manifestation of a mixing condition that provides a certain independence of the long-term behaviour of the process with respect to any finite initial block/window.

In Proposition 1 below we demonstrate that processes with *i.i.d.* increments naturally satisfy Assumptions 1-3. Before proving this we define a coercive function as follows.

*Definition 4:* A function  $f : \mathbb{R}_+ \mapsto \mathbb{R}_+ \cup \{+\infty\}$  with domain  $\text{Dom}(f) := \{x : f(x) < +\infty\}$  is defined to be *coercive* if for every  $y \in \mathbb{R}_+$ , there is compact set  $\mathcal{C} \subsetneq \text{Dom}(f)$  such that  $f(x) > y$  for all  $x \in \text{Dom}(f) \setminus \mathcal{C}$ .

*Proposition 1:* For scale parameter  $L$  assume that the arrival process has *i.i.d.* increments while the arrival processes for different users are independent. Furthermore, assume that for every  $k$ ,  $\{A_1^{k,L}\}$  satisfies the conditions of the Gärtner-Ellis Theorem [51, Thm. 2.3.6] and that the limiting rate function is either coercive or has domain  $\mathfrak{R}_+$ . Then Assumptions 1-3 hold.

*Proof:* The fact the above class of processes satisfy Assumption 1 is immediate from [24, Thm. 3]. Similarly the Gärtner-Ellis Theorem [51, Thm. 2.3.6] yields convexity and from the coercivity of the rate function or with its domain being  $\mathbb{R}_+$ , Assumption 2 follows. When the arrival process has *i.i.d.* increments, then it also follows that for  $\mathbf{y} \in \mathbb{R}_+^{K \times t}$  and  $\mathbf{x} \in \mathcal{D}^K$ ,

$$I_t^\sharp(\mathbf{y}) = \sum_{k=1}^K \sum_{s=1}^t \Lambda_1^{*,k}(y_s^k),$$

and

$$I^\sharp(\mathbf{x}) = \sum_{k=1}^K \sum_{s=1}^{+\infty} \Lambda_1^{*,k}(x_s^k),$$

where  $\Lambda_1^{*,k}$  is the Fenchel-Legendre transform of  $\Lambda_1^k(\theta) := \lim_{L \rightarrow \infty} \frac{1}{L} \log E e^{\theta L A_1^{k,L}}$ . Note that that  $I^\sharp(\mathbf{x}) = +\infty$  if  $\mathbf{x} \in \mathcal{D}^K \setminus \mathcal{D}_\mu^K$ .

Now for every sequence  $\{\mathbf{x}^n\}$  converging to  $\mathbf{x}$  in  $\mathcal{D}_\mu^K$  with  $I^\sharp(\mathbf{x}) < +\infty$  such that there exists  $t \in \mathbb{N}$  so that  $\mathbf{x}_s^n = \mathbf{x}_s$  for all  $s > t$  (for all  $n$ ), we have  $I^\sharp(\mathbf{x}^n) < +\infty$  for all  $n$  large enough

and

$$I^\sharp(\mathbf{x}^n) - I^\sharp(\mathbf{x}) = I^\sharp_t(\mathbf{x}^n|_{(0,t]}) - I^\sharp_t(\mathbf{x}|_{(0,t]}) \xrightarrow{n \rightarrow +\infty} 0.$$

■

*Remark 6:* A more general characterization, especially of Assumption 3, remains an important area of future work. However, to provide insight about Assumptions 2-3 and their relationship to Assumption 1 we provide the following simple example. Our example will be constructed as averages of  $L$  stationary *i.i.d.* random processes so it is sufficient to describe the underlying stochastic process: We let the random process be  $X_k = \gamma + \gamma V$  when  $k$  is odd and  $X_k = \gamma - \gamma V$  when  $k$  is even, where  $V$  is a Uniform $[-1, 1]$  random variable and  $\gamma > 0$  is the long-term mean for all our sequences. Here it can be verified that  $I(x) < +\infty$  if and only if  $x_k = \gamma(1 - \alpha)$  if  $k$  is even and  $x_k = \gamma(1 + \alpha)$  for  $k$  odd where  $\alpha \in [-1, 1]$ . It can also be verified that  $I(x) = 0$  for the all- $\gamma$  sequence, i.e.,  $x_k \equiv \gamma$ . However, it is clear that one can easily construct sequences as required in Assumption 3 with the rate function being infinite which, nevertheless, converge to the all- $\gamma$  sequence. In other words, Assumptions 2-3 require more than the regularity conditions from [24], [26] that guarantee validity of Assumption 1.

Given the above assumptions on the arrival processes, we are now ready to provide an overview of the main results of the paper.

### B. Main Results: An Overview

Assuming that the sequence of the arrival processes  $\{\mathbf{A}^L\}$  satisfies a many-sources sample-path LDP with a “continuous” rate function (Assumptions 1-3), LDPs for the finite and infinite-horizon workloads will be a direct consequence of Garcia’s extended contraction principle once the required quasi-continuity and almost compactness properties are shown. We will demonstrate that the quasi-continuity and almost compactness of the finite horizon workload mappings are inherited from the quasi-continuity of the schedulers  $H$  and  $H^{\text{wc}}$ .

*Theorem 1:* Under Assumptions 1-2 and for all  $t \in \mathbb{N}$ , the sequence of the finite-horizon workloads under a max-weight scheduler  $\{\mathbf{W}_{0,t}^L := G_t(\mathbf{A}^L|_{(0,t]})\}$  satisfies an LDP on  $\mathbb{R}_+^K$  with the rate function  $I_t$ , where for  $\mathbf{b} \in \mathbb{R}_+^K$

$$I_t(\mathbf{b}) = \inf_{\mathbf{x} \in \mathbb{R}_+^{K \times t}: \mathbf{x} \mathbf{G}_t \ni \mathbf{b}} I^\sharp_t(\mathbf{x}). \quad (7)$$

*Theorem 2:* Under Assumptions 1-2 and for all  $t \in \mathbb{N}$ , the sequence of the finite-horizon workloads under the work-conserving max-weight scheduler,  $\{\mathbf{W}_{0,t}^L := G_t^{\text{wc}}(\mathbf{A}^L|_{(0,t]})\}$ , satisfies an LDP on  $\mathbb{R}_+^K$  with the rate function  $I_t$ , where for  $\mathbf{b} \in \mathbb{R}_+^K$

$$I_t(\mathbf{b}) = \inf_{\mathbf{x} \in \mathbb{R}_+^{K \times t} : \mathbf{x} \mathbf{G}_t^{\text{wc}} \ni \mathbf{b}} I_t^\#(\mathbf{x}). \quad (8)$$

The quasi-continuity and almost compactness of the stationary workload mapping, however, requires slightly more work as shown in Section VI. In fact, unlike the finite-horizon LDP results of Theorems 1 and 2, which hold for general rate regions and for both max-weight and its work conserving version, the infinite horizon LDP result of Theorem 3 is established only under the work conserving max-weight policy and only with a simplex rate region. For a vector  $\mathbf{x} \in \mathbb{R}_+^K$  define  $\hat{\mathbf{x}} := \sum_{k=1}^K x^k / C^k$ , then the simplex rate region is given by,

$$\mathcal{R}_s := \{\mathbf{r} \in \mathbb{R}_+^K : \hat{\mathbf{r}} \leq 1\}. \quad (9)$$

*Theorem 3:* Consider a work conserving max-weight scheduler with a simplex rate region  $\mathcal{R}_s$ . Let  $\boldsymbol{\mu} \in \mathbb{R}^K$  be an admissible arrival rate vector strictly inside this simplex rate region, i.e.,  $\boldsymbol{\mu} \in \text{int}(\mathcal{R}_s)$  and  $\{\mathbf{A}^L\}$  be a sequence of arrival processes that satisfies Assumptions 1-3 in  $\mathcal{D}_\mu^K$ . The sequence of infinite-horizon workloads  $\{\mathcal{W}^L := G^{\text{wc}}(\mathbf{A}^L)\}$  satisfies an LDP on  $\mathbb{R}_+^K$  with good rate function  $J$ , where for  $\mathbf{b} \in \mathbb{R}_+^K$

$$J(\mathbf{b}) = \inf_{\mathbf{a} \in \mathcal{D}_\mu^K : \mathbf{a} \mathbf{G}^{\text{wc}} \ni \mathbf{b}} I^\#(\mathbf{a}). \quad (10)$$

Note that all the rate functions have the appearance of being what one would naturally expect, i.e., among all arrival sequences that result in the required workload at the required epoch, find the one with the least cost to deduce the rate function. However, the discontinuity of the queueing map makes this a non-trivial assertion and also enlarges the set of allowed arrival sequences.

## VI. ANALYSIS: LDP'S FOR WORKLOADS

In this section, we prove the main results of the paper: establishing LDP for the sequences of the finite-horizon and infinite-horizon workloads. We first delineate the proof for the LDPs for the sequence of the finite-horizon workloads.

### A. LDP for Finite-Horizon Workloads

In this section, for  $t \in \mathbb{N}$ , we establish an LDP for finite-horizon workloads  $\{\mathbf{W}_{0,t}^L := G_t^{\text{wc}}(\mathbf{A}^L|_{(0,t]})\}$  and  $\{\mathbf{W}_{0,t}^L := G_t(\mathbf{A}^L|_{(0,t]})\}$ . The approach is to first show that the mappings  $G_t^{\text{wc}} : \mathbb{R}_+^{Kt} \mapsto \mathbb{R}_+^K$  and  $G_t : \mathbb{R}_+^{Kt} \mapsto \mathbb{R}_+^K$  are quasi-continuous and almost compact and use Garcia's extended contraction principle to obtain an LDP for the finite-horizon workloads from the LDP assumption for  $\{\mathbf{A}^L|_{(0,t]}\}$ . From Fact 1, the almost compactness and quasi-continuity of workload mappings are sufficient to establish an LDP for finite-horizon workload, since according to Assumption 2,  $I_t^\sharp(\cdot)$  is continuous (in the product topology).

Using the quasi-continuity of the scheduling function we now prove the required quasi-continuity of the workload maps. First we consider the work-conserving max-weight scheduler.

*Lemma 2:* For  $t \in \mathbb{N}$ ,  $G_t^{\text{wc}}$  is quasi-continuous with respect to the scaled uniform topology.

*Proof:* See Appendix B. The induction-based proof uses the quasi-continuity of the scheduler  $H$  and the linear dependence of the workload  $\mathbf{W}_s$  at time  $-s$  on  $\mathbf{A}_{s+1}$  for all  $s \in (0, t-1]$ . ■

Similarly, in order to establish an LDP for the workload process under the maximum weight scheduler, we have the quasi-continuity of function  $G_t$ .

*Lemma 3:* For  $t \in \mathbb{N}$ ,  $G_t(\cdot)$  is quasi-continuous with respect to the scaled uniform topology.

*Proof:* See Appendix B. In this induction-based proof we establish certain analytical properties of the workload map for all possible rate-regions that then allows us to obtain a quasi-continuous selection. ■

Having established quasi-continuity, the next result demonstrates the almost compactness of the workload maps.

*Lemma 4:* For  $t \in \mathbb{N}$ , both  $G_t$  and  $G_t^{\text{wc}}$  are almost compact on  $\mathbb{R}_+^K$  with respect to the scaled uniform topology.

*Proof:* Since the workload at any time cannot exceed the amount work brought in from the last time the system was empty irrespective of the scheduler used, we automatically get the following bounds

$$G_t(\mathbf{a}|_{(0,t]}) \leq \mathbf{a}(0, t], \quad G_t^{\text{wc}}(\mathbf{a}|_{(0,t]}) \leq \mathbf{a}(0, t].$$

This implies the almost compactness of  $G_t$  and  $G_t^{\text{wc}}$ . Since every sequence  $\{\mathbf{a}^n\}$  converging to  $\mathbf{a}$  (in  $\mathcal{D}_\mu^K$  in the scaled uniform norm topology) converges point-wise too, it follows that

$G_t(\mathbf{a}^n|_{(0,t]})$  and  $G_t^{\text{wc}}(\mathbf{a}^n|_{(0,t]})$  will lie in a bounded subset of  $\mathbb{R}_+^K$  so we can use the Bolzano-Weierstrass theorem. ■

Now, as already discussed, the proof of Theorems 1 and 2 is complete.

Next, we discuss the LDP for the infinite-horizon workloads for the work-conserving max-weight scheduler with a simplex rate-region.

### B. LDP for Infinite-Horizon Workloads

In this section, we establish an LDP of the sequence of the infinite-horizon workloads  $\{\mathcal{W}^L = G^{\text{wc}}(\mathbf{A}^L)\}$  where  $\mathbf{A}^L \in \mathcal{D}_\mu^K$ ,  $\mu \in \text{int}(\mathcal{R}_s)$ , and  $\{\mathbf{A}^L\}$  satisfies Assumptions 1-3. Similar to the last section, we first show that the mapping  $G^{\text{wc}}$  is quasi-continuous and almost compact on  $\mathcal{D}_\mu^K$  and then use Garcia's extended contraction principle to establish the desired LDP.

The following lemmas establish the necessary steps to apply Garcia's contraction principle to the stationary workload map. The first of these lemmas relates the infinite horizon workload mapping to that of a finite horizon.

*Lemma 5:* Consider an arrival process  $\mathbf{a} \in \mathcal{D}_\mu^K$ . There exists a  $s^* = s^*(\mathbf{a}) < \infty$  such that the workloads at time  $-s^*$  under  $\mathbf{a}$  falls within the rate region  $\mathcal{R}_s$ , i.e.,  $G^{\text{wc}}(\mathbf{a}|_{(s^*,\infty)}) \in \mathcal{R}_s$ . Furthermore, for any sequence of arrival processes  $\{\mathbf{a}^n \in \mathcal{D}_\mu^K\}$  converging to  $\mathbf{a}$  (in scaled uniform topology), the workloads at time  $-s^*$  under  $\mathbf{a}^n$ , when  $n$  is large enough, also fall within the rate region  $\mathcal{R}_s$ , i.e.,  $\exists n_0$  such that  $G^{\text{wc},\mathcal{R}_s}(\mathbf{a}^n|_{(s^*,\infty)}) \in \mathcal{R}_s$  for  $n > n_0$ .

*Proof:* See Appendix C. The main idea is to use the fact that a suitably normalized sum (over all queues) workload process behaves like the workload of a single server queue with only one flow and a work-conserving service discipline. This allows us to use the continuity of the workload mapping of a single server queue and the stability of the queue to arrive at the assertion of the lemma. ■

Now we are ready to verify the requirements of Garcia's extended contraction principle.

*Lemma 6:* Let  $\mathbf{a} \in \mathcal{D}_\mu^K$  be an arrival process with rate  $\mu \in \text{int}(\mathcal{R}_s)$ , and  $G^{\text{wc}}(\mathbf{a})$  be the corresponding infinite-horizon workload. For any  $\mathcal{W} \in {}^{\mathbf{a}}G^{\text{wc}}$  there exists a sequence of arrivals  $\{\mathbf{a}^n \in \mathcal{D}_\mu^K\}$  such that  $\mathbf{a}^n$  converges to  $\mathbf{a}$  in scaled uniform topology,  $G^{\text{wc}}(\mathbf{a}^n) \rightarrow \mathcal{W}$ ,  $G^{\text{wc}}$  is continuous at  $\mathbf{a}^n$ , and  $I^\sharp(\mathbf{a}^n) \rightarrow I^\sharp(\mathbf{a})$ .

*Proof:* See Appendix C. The idea is to relate the infinite horizon workload, using Lemma 5, to a finite-horizon workload mapping whose quasi-continuity was established earlier. ■

*Lemma 7:* The mapping  $G^{\text{wc}, \mathcal{R}_s}$  is almost compact on  $\mathcal{D}_\mu^K$  with respect to the scaled uniform topology.

*Proof:* See Appendix. The proof is along the same lines as that of Lemma 6. ■

Again, the above lemmas and Garcia's extended contraction principle immediately yield the LDP for the sequence of the infinite-horizon workload in Theorem 3.

Let us now consider the problem of calculating the rate function. Eqn. (10) suggests that the rate function  $J$ , where  $J(\mathbf{b}) = \inf_{\mathbf{a} \in \mathcal{D}_\mu^K: \mathbf{a} \mathbf{G}^{\text{wc}} \ni \mathbf{b}} I^\sharp(\mathbf{a})$ , could be interpreted as the minimum-cost solution among all paths  $\mathbf{a} \in \mathcal{D}_\mu^K$  such that  $\mathbf{b} \in \mathbf{a} \mathbf{G}^{\text{wc}}$ , where the cost of the path  $\mathbf{a}$  is  $I^\sharp(\mathbf{a})$  and convex. Hence, the problem of finding the rate functions is a deterministic optimal control problem like those in [12], [14]. However, the expressions for the rate functions  $I_t$  and  $J$  in (7) and (10) are of little use in their current forms, as their computation is far from straight forward. In the next section, we simplify the rate functions when the arrival processes are limited to having *i.i.d.* increments.

## VII. I.I.D. INCREMENTS: SIMPLIFIED RATE FUNCTIONS

In this section, we give a calculation of the finite-horizon and infinite-horizon rate functions in the case when the arrivals have *i.i.d.* increments. In this case, the cost of a sample path  $\mathbf{a} \in \mathcal{D}^K$ , which is  $I^\sharp(\mathbf{a})$ , is additive and the total cost of any arrival sample path is the sum of the cost over all timeslots and queues. This property helps us to simplify the calculation of the rate functions.

Consider the underlying arrival process  $\mathbf{A}$  to be a process with *i.i.d.* increments, e.g., a compound Poisson arrival process with exponential packet length (see [40]). For these *i.i.d.* increment arrival processes, as mentioned earlier, we have

$$I_t^\sharp(\mathbf{a}|_{(0,t]}) = \sum_{k=1}^K \sum_{i=1}^t \Lambda_1^*(x_i^k), \text{ and } I^\sharp(\mathbf{a}) = \sum_{k=1}^K \sum_{i=1}^{+\infty} \Lambda_1^*(x_i^k).$$

Next, we simplify these rate functions.

### A. Infinite-Horizon Rate Function

The following lemma expresses the infinite-horizon rate function  $J$  as the infimum of the finite-horizon rate functions  $I_t$  over all time  $t$ .

*Lemma 8:* For *i.i.d.* increment arrival processes with  $\boldsymbol{\mu} \in \text{int}(\mathcal{R}_s)$ , the infinite-horizon rate function  $J$  is simplified as

$$J(\mathbf{b}) = \inf_{t \geq 1} I_t(\mathbf{b}). \quad (11)$$

*Proof:* The cost of a sample path over time is the sum of the cost of arrivals in all timeslots. As in the proof of Lemma 6, for  $\mathbf{a} \in \mathcal{D}_\mu^K$  where  $\boldsymbol{\mu} \in \text{int}(\mathcal{R}_s)$ , we can find  $t := s^*(\mathbf{a})$  such that  $\mathbf{W}_t(\mathbf{a}) \in \mathcal{R}_s$ . Hence, for  $\mathbf{a}$  such that  $\mathbf{b} \in {}^{\mathbf{a}}\mathbf{G}^{\text{wc}}$ , one can reduce the cost of the path by constructing a new sample-path  $\tilde{\mathbf{a}}$  by setting  $\tilde{\mathbf{a}}_v = \boldsymbol{\mu}$  for all  $v > t$  and  $\tilde{\mathbf{a}}_v = \mathbf{a}_v$  for all  $v \leq t$  while still satisfying  $\mathbf{b} \in \tilde{\mathbf{a}}\mathbf{G}^{\text{wc}}$ . This is because  $I^\sharp(\tilde{\mathbf{a}}) = I^\sharp_t(\mathbf{a}|_{(0,t]}) \leq I^\sharp(\mathbf{a})$ . On the other hand, since  $\mathbf{W}_t(\mathbf{a}) \in \mathcal{R}_s$ , we can write  $\mathbf{b} \in {}^{\mathbf{a}|_{(0,t]}}\mathbf{G}_t^{\text{wc}}$ . All of these imply that

$$J(\mathbf{b}) = \inf_{\mathbf{a} \in \mathcal{D}_\mu^K: {}^{\mathbf{a}}\mathbf{G}^{\text{wc}} \ni \mathbf{b}} I^\sharp(\mathbf{a}) = \inf_{t \geq 1} \inf_{\mathbf{x} \in \mathbb{R}_+^{Kt}: {}^{\mathbf{x}}\mathbf{G}_t^{\text{wc}} \ni \mathbf{b}} I^\sharp_t(\mathbf{x}) = \inf_{t \geq 1} I_t(\mathbf{b}),$$

by the definition of  $I_t(\mathbf{b})$  in (7). ■

With this simplification available, we now look at the finite-horizon rate function  $I_t$  in more details.

### B. Finite-Horizon Rate Function

In this subsection, we provide a further simplified expression of the finite-horizon rate function  $I_t$ .

*Lemma 9:* For  $t \in \mathbb{N}$ , the finite-horizon rate function  $I_t$  is simplified as

$$I_t(\mathbf{b}) = \min \left( I^\sharp_1(\mathbf{b}), \min_{u \in (1,t]} \inf_{\mathbf{x} \in \mathbb{A}(u,\mathbf{b})} I^\sharp_u(\mathbf{x}) \right) \quad (12)$$

for  $\mathbf{b} \in \mathbb{R}_+^K$ , where for  $u > 1$

$$\mathbb{A}(u, \mathbf{b}) := \{\mathbf{x} \in \mathbb{R}_+^{K \times u} : \mathbf{b} \in {}^{\mathbf{x}}\mathbf{G}_u^{\text{wc}}, G_{u-v}^{\text{wc}}(\mathbf{x}|_{(v,u]}) \notin \mathcal{R}_s, \forall v \in [1, u-1]\}. \quad (13)$$

*Proof:* This follows the idea from the proof of Lemma 8. Let  $t \in \mathbb{N}$ . For  $\mathbf{a} \in \mathcal{D}_\mu^K$  such that  $\mathbf{b} \in {}^{\mathbf{a}}\mathbf{G}_t^{\text{wc}}$ , we let

$$u = \min \left\{ t, \min \left\{ s \in [1, t-1] : \mathbf{W}_s = G_{t-s}^{\text{wc}}(\mathbf{a}|_{(s,t]}) \in \mathcal{R}_s \right\} \right\}.$$

In other words,  $-u$  is the last time the workload vector is inside the capacity region  $\mathcal{R}_s$  before time 0. With this definition of  $u$ , for all  $u > 1$  we have  $\mathbf{W}_v \notin \mathcal{R}_s$  for all  $v \in [1, u-1]$ . By definition of  $I_t$ , we already know that the workload vector starts initially inside  $\mathcal{R}_s$  at time  $-t$ .

Therefore, we can find another path  $\tilde{\mathbf{a}} \in \mathbb{R}_+^{Kt}$  with a reduced cost while keeping the workloads at time  $-u+1$  to 0 (i.e.,  $\mathbf{W}_{u-1}$  to  $\mathbf{W}_0$ ) intact by setting  $\tilde{\mathbf{a}}_v = \boldsymbol{\mu}$ ,  $\forall v \in (u, t]$  and  $\tilde{\mathbf{a}}_v = \mathbf{a}_v$  otherwise. It is easy to see that we have  $I_t^\sharp(\mathbf{a}|_{(0,t]}) \geq I_u^\sharp(\mathbf{a}|_{(0,u]}) = I_t^\sharp(\tilde{\mathbf{a}}|_{(0,t]})$  and yet  $\mathbf{b} \in \tilde{\mathbf{a}}|_{(0,u]} G_u^{\text{wc}} = \tilde{\mathbf{a}} G_t^{\text{wc}}$ . The same logic also applies to the case when  $u = 1$ . However, in this case the only way that we can achieve a workload vector  $\mathbf{b}$  at 0 is if exactly that amount of work arrives, i.e., if  $\mathbf{a}_1 = \mathbf{b}$ .

Since by definition  $\mathbf{W}_v = G_{u-v}^{\text{wc}}(\mathbf{a}|_{(v,u]})$  for  $v \in [1, u-1]$  (when  $u > 1$ ), we have

$$\begin{aligned} I_t(\mathbf{b}) &= \inf_{\mathbf{x} \in \mathbb{R}_+^{Kt}: \mathbf{x} G_t^{\text{wc}} \ni \mathbf{b}} I_t^\sharp(\mathbf{x}|_{(0,t]}) = \min \left( I_1^\sharp(\mathbf{b}), \min_{u \in (1,t]} \inf_{\substack{\mathbf{x} \in \mathbb{R}_+^{Ku}: \mathbf{b} \in \mathbf{x} G_u^{\text{wc}}, \\ G_{u-v}^{\text{wc}}(\mathbf{x}|_{(v,u]}) \notin \mathcal{R}_s}} I_u^\sharp(\mathbf{x}) \right) \\ &= \min \left( I_1^\sharp(\mathbf{b}), \min_{u \in (1,t]} \inf_{\mathbf{x} \in \mathbb{A}(u, \mathbf{b})} I_u^\sharp(\mathbf{x}|_{(0,u]}) \right), \end{aligned}$$

where  $\mathbb{A}(u, \mathbf{b})$  is defined as in (13). ■

*Remark 7:* The above lemma reduces the set of feasible sample paths to the set  $\mathbb{A}(u, \mathbf{b})$  for  $u \in (0, t]$ . It is interesting to note the property of the sample-paths in this set. For any  $\mathbf{x} \in \mathbb{A}(u, \mathbf{b})$ , we have  $\hat{\mathbf{W}}_0(\mathbf{x}) = \hat{\mathbf{x}}(0, u] - (u-1) = \hat{\mathbf{b}}$ , recalling that the  $\hat{\cdot}$  notation is the normalized sum of the elements of the vectors. There is no wastage of service capacity over the  $u-1$  timeslots because  $\forall v \in [1, u-1]$ ,  $\mathbf{W}_v = G_{u-v}^{\text{wc}}(\mathbf{x}|_{(v,u]}) \notin \mathcal{R}_s$  and hence  $\hat{\mathbf{W}}_v > 1$ . That is, any sample path  $\mathbf{x} \in \mathbb{A}(u, \mathbf{b})$  has its normalized sum of the arrivals over time  $(0, u]$  and queues equal to  $\hat{\mathbf{x}}(0, u] = \hat{\mathbf{b}} + (u-1)$ .

In addition, an immediate implication of Lemma 9 is that we can rewrite  $J$  in (10) as

$$\begin{aligned} J(\mathbf{b}) &= \inf_{t \geq 1} I_t(\mathbf{b}) = \inf \left( I_1(\mathbf{b}), \inf_{t \geq 2} I_t(\mathbf{b}) \right) = \inf \left( I_1^\sharp(\mathbf{b}), \inf_{t \geq 2} \min_{u \in (1,t]} \inf_{\mathbf{x} \in \mathbb{A}(u, \mathbf{b})} I_u^\sharp(\mathbf{x}|_{(0,u]}) \right) \\ &= \inf \left( I_1^\sharp(\mathbf{b}), \inf_{t \geq 2} \inf_{\mathbf{x} \in \mathbb{A}(t, \mathbf{b})} I_t^\sharp(\mathbf{x}|_{(0,t]}) \right), \quad (14) \end{aligned}$$

where we also used the fact that  $I_1(\mathbf{b}) = I_1^\sharp(\mathbf{b})$ . If we denote  $t^*$  as the optimizer of the last equation, then  $t^*$  is called the *critical timescale* (see [24]). It can then be interpreted that  $t^*$  is the most likely length of time it would take to “fill” the buffers to a given level  $\mathbf{b}$  from being “empty” (more precisely, anywhere within  $\mathcal{R}_s$ ).

*Remark 8:* The induction-based proof of Lemma 2 reveals another important property of the sample-paths, namely, that at every stage it suffices to consider the quasi-continuous selection  $H^{\text{wc}}(\cdot)$  of the scheduling function. This then implies that we need to consider all valid arrivals sequences that respect the constraints of the set  $\mathbb{A}(\cdot, \mathbf{b})$  such that using any of the allowed (based

on the workload vector at each time) scheduling actions given by  $H^{\text{wc}}(\cdot)$  results in the workload vector  $\mathbf{b}$  at the required epoch.

Note that for fixed  $u \in \mathbb{N}$ ,  $\inf_{\mathbf{x} \in \mathbb{A}(u, \mathbf{b})} I_u^\#(\mathbf{x}|_{(0, u]})$  is an optimization problem, with a convex cost function  $I_u^\#(\cdot)$  and a set  $\mathbb{A}(u, \mathbf{b})$  of allowable solutions, in general, a  $K(u-1)$ -dimensional set. This problem is difficult to solve analytically. Since the cost function  $I_u^\#(\mathbf{x}|_{(0, u]}) = \sum_{i=1}^u \sum_{k=1}^K \Lambda_1^{*,k}(\mathbf{x}_i)$  is additive, a possible numerical method is the numerical backwards induction of dynamic programming. However, the method suffers from the curse of dimensionality and hence is not practical for large  $u$  and  $\mathbf{b}$ . Hence, we turn our attention to finding some simplified bounds of the rate functions. This can be done by employing the additivity and convexity of the rate function  $I_t^\#$ . Next we present some bounds for the case when  $K = 2$ .

### C. Properties of the Minimum-Cost Sample Paths

Here, we see that the convexity of the cost function  $\Lambda_1^{*,k}$  for all  $k \in \mathcal{K}$ , induces two properties for the optimal paths.

*Property 1: Constant-speed linear path is the cheapest.* Among all arrival sample paths  $\mathbf{x} \in \mathcal{D}_\mu^K$  with the only constraint that  $\mathbf{x}(0, t] = \mathbf{b}$ , i.e., the total amount of arrivals at the end of time  $t \in \mathbb{N}$  equals  $\mathbf{b}$ , the cheapest or minimum-cost path is the constant-speed linear path, where the arrival in each timeslot is equal to  $\mathbf{b}/t$ .

*Proof:* This is because the path cost function is additive, i.e.,  $\Lambda_t^*(\mathbf{x}) = \sum_{k=1}^K \sum_{i=1}^t \Lambda_1^{*,k}(x_i^k)$ , and the per-timeslot cost function  $\Lambda_1^{*,k}$  is convex. Applying Jensen's inequality [54] gives

$$\Lambda_t^*(\mathbf{x}) = \sum_{k=1}^K \sum_{i=1}^t \Lambda_1^{*,k}(x_i^k) \geq \sum_{k=1}^K t \Lambda_1^{*,k} \left( \frac{1}{t} \sum_{i=1}^t x_i^k \right) = \sum_{k=1}^K t \Lambda_1^{*,k}(b^k/t),$$

with equality when  $x_i^k = b^k/t$  for all  $i$  and  $k$ . See an illustration in Figure 1(a). ■

From now onwards, without loss of generality, we assume that the arrivals, workloads and service vectors are scaled by  $1/\mathbf{C}$  where  $\mathbf{C} = (C^1, \dots, C^K)$  is the vector of maximum service rates. In this normalized setting, the maximum service rate of all the users is 1. Assume also that  $\Lambda_1^{*,k}$  for  $k \in \mathcal{K}$  is suitably modified for this scaling, and, for ease of exposition, that with this scaling the processes are statistically identical, i.e.,  $\Lambda_1^{*,k} \equiv \Lambda_1^*$  for all  $k \in \mathcal{K}$ . We can then write down the following property.

*Property 2: Constant-speed linear path closest to the equal line is the cheapest.* For constant-speed linear paths  $\mathbf{a} \in \mathbb{R}_+^{Kt}$  with the sum  $\hat{\mathbf{a}}(0, t] = d$ , the cost of the path is a Schur-convex

function [55].

*Proof:* Since the arrival paths are constant-speed linear path, without loss of generality we can consider arrival paths in a single timeslot. Consider path  $\mathbf{x} \in \mathbb{R}_+^K$ , then the cost of the path is  $\sum_{k=1}^K \Lambda_1^*(x^k)$  where  $\Lambda_1^*(\cdot)$  is a convex function. Therefore from the results in [55] it follows that  $\sum_{k=1}^K \Lambda_1^*(x^k)$  is a Schur-convex function. In order words, if  $\mathbf{x}$  is majorized by  $\mathbf{y} \in \mathbb{R}_+^K$  (denoted by  $\mathbf{x} \prec \mathbf{y}$ ), i.e.,  $\sum_{k=1}^j x^{[i]} \leq \sum_{k=1}^j y^{[i]}$  for  $j = 1, \dots, K-1$  and  $\sum_{k=1}^K x^k = \sum_{k=1}^K y^k$  where  $x^{[i]}$  is the  $i^{\text{th}}$  largest component of  $\mathbf{x}$ , then  $\sum_{k=1}^K \Lambda_1^*(x^k) \leq \sum_{k=1}^K \Lambda_1^*(y^k)$ .

This is easily appreciated when  $K = 2$ . Let  $\mathbf{x} = (x, d-x) \in \mathbb{R}_+^2$  and  $\mathbf{y} = (y, d-y) \in \mathbb{R}_+^2$ , where  $y > x > d/2$ , then we have  $\Lambda_1^*(x) + \Lambda_1^*(d-x) \leq \Lambda_1^*(y) + \Lambda_1^*(d-y)$ , and  $\mathbf{x}$  is cheaper than  $\mathbf{y}$ . We illustrate this in Figure 1(b). ■

These properties are also used in [9], [12], [13] for large-deviations analysis of scheduling disciplines. Next, we use these properties to calculate  $I_2$  and bounds on  $I_t$  for  $t \in \mathbb{N}$  for just the work-conserving scheduler operating on the simplex rate-region.

#### D. Example: Calculation of $I_2$

Here we look at an example for calculation of the finite-horizon rate function  $I_t$  to illustrate that the calculation continues to be rather involved. For simplicity, consider the case when  $t = 2$  and  $K = 2$ . From (12),  $I_2(\mathbf{b})$  for  $\mathbf{b} \in \mathbb{R}_+^2$  can be written as

$$I_2(\mathbf{b}) = \min \left\{ \sum_{i=1}^2 \Lambda_1^*(b^i), \inf_{\mathbf{x}|_{(0,2]} \in \mathbb{A}(2,\mathbf{b})} \sum_{i=1}^2 \sum_{k=1}^2 \Lambda_1^*(x_i^k) \right\}, \quad (15)$$

where

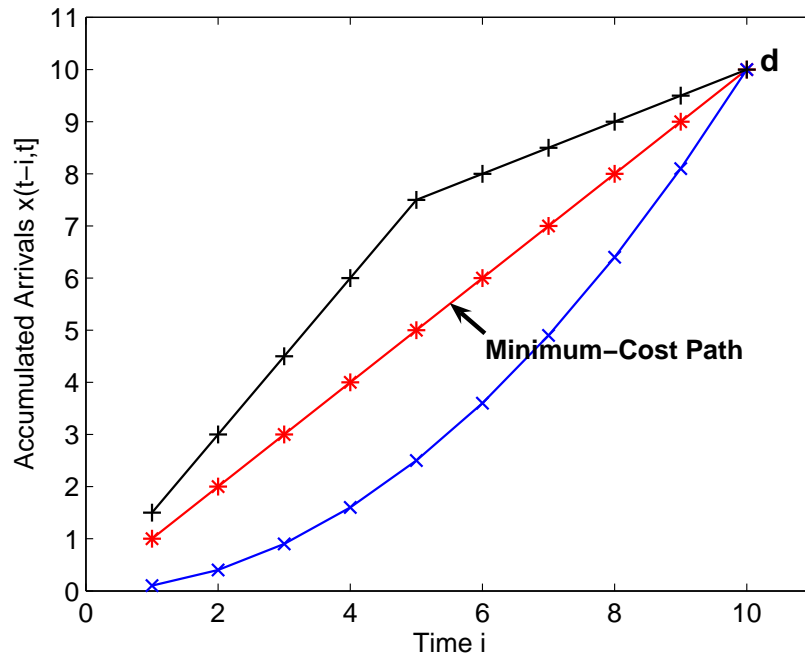
$$\mathbb{A}(2, \mathbf{b}) = \left\{ \mathbf{a}|_{(0,2]} \in \mathbb{R}_+^4 : \mathbf{a}_2 \notin \mathcal{R}_s, \mathbf{b} \in \mathbf{a}|_{(0,2]} G_2^{\text{wc}} \right\}.$$

The workload at time zero is  $\mathbf{W}_0 = G_2^{\text{wc}}(\mathbf{a}|_{(0,2]}) = \mathbf{a}_1 + [\mathbf{a}_2 - H^{\text{wc}}(\mathbf{a}_2)]^+$ , which is equal to  $\mathbf{a}(0, 2] - H^{\text{wc}}(\mathbf{a}_2)$  since  $\mathbf{a}_2 \notin \mathcal{R}_s$ . On the other hand, we require  $\mathbf{W}_0 = \mathbf{b}$ . Hence, using the scheduler  $H^{\text{wc}}$  given in (3), we can express  $\mathbb{A}(2, \mathbf{b})$  as  $\mathbb{A}(2, \mathbf{b}) = \mathbb{A}_{(1)} \cup \mathbb{A}_{(2)} \cup \mathbb{A}_{(3)}$ , where  $\mathbb{A}_{(j)} \subseteq \mathbb{R}_+^4, j = 1, 2, 3$ , are defined as

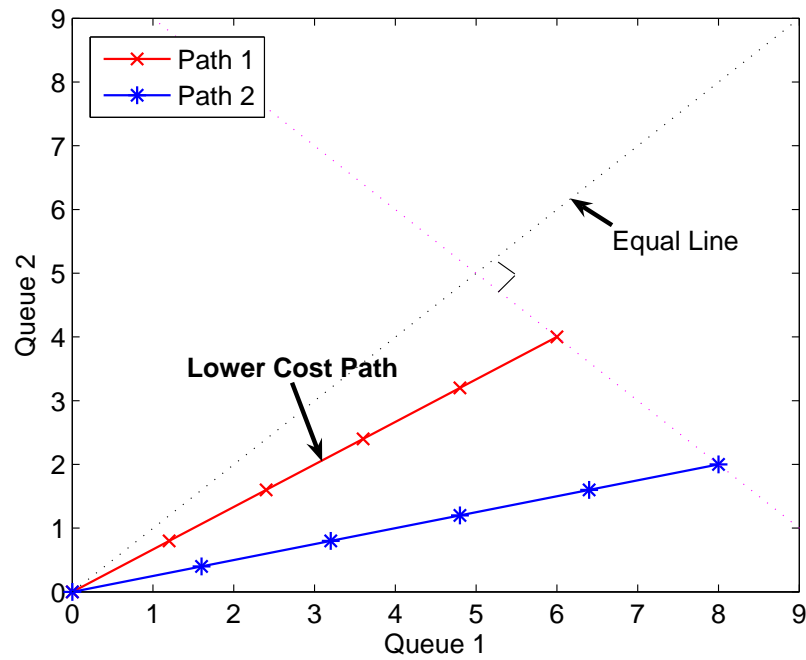
$$\mathbb{A}_{(1)} := \{(\mathbf{a}_1, \mathbf{a}_2) \in \mathbb{R}_+^4 : a_2^1 \geq a_2^2, a_2^1 \geq 1, \mathbf{a}(0, 2] = \mathbf{b} + (1, 0)\},$$

$$\mathbb{A}_{(2)} := \{(\mathbf{a}_1, \mathbf{a}_2) \in \mathbb{R}_+^4 : a_2^1 \leq a_2^2, a_2^2 \geq 1, \mathbf{a}(0, 2] = \mathbf{b} + (0, 1)\},$$

$$\mathbb{A}_{(3)} := \{(\mathbf{a}_1, \mathbf{a}_2) \in \mathbb{R}_+^4 : a_2^1 \leq 1, a_2^2 \leq 1, a_2^1 + a_2^2 \geq 1, \mathbf{a}(0, 2] = \mathbf{b} + \text{Proj}_{\mathcal{R}_s}(\mathbf{a}_2)\}.$$



(a) Property 1



(b) Property 2

Fig. 1. Two properties of the minimum-cost sample paths: (a) Property 1: the minimum-cost path is the constant-speed linear path. (b) Property 2: Path 1 which is closer to the Equal Line has a lower cost than Path 2.

Note that in the definition of  $\mathbb{A}_{(1)}$  and  $\mathbb{A}_{(2)}$  we have used the property highlighted in Remark 8. For the two user case the scheduling function  $H^{\text{wc}}(\cdot)$  is only discontinuous at  $\mathbf{x} \in \mathbb{R}_+^2$  such that  $x^1 = x^2 \geq 1$ . Here one can either choose the service vector  $(1, 0)$  or  $(0, 1)$  to obtain a quasi-continuous selection. Thus, both of these options have to be considered.

Using these newly defined sets the second term in the RHS of (15) can be rewritten as

$$\inf_{\mathbf{x} \in \mathbb{A}(2, \mathbf{b})} \sum_{i=1}^2 \sum_{k=1}^2 \Lambda_1^*(x_i^k) = \min_{j \in [1, 3]} \inf_{\mathbf{x} \in \mathbb{A}_{(j)}} \sum_{i=1}^2 \sum_{k=1}^2 \Lambda_1^*(x_i^k).$$

Trajectories of some examples of the (accumulated) arrival sample paths are illustrated in Figure 2(a) and their corresponding workload trajectories in Figure 2(b). In particular, Figure 2(a) shows example trajectories of the accumulated arrival sample paths  $\mathbf{a}|_{(0,2]} \in \mathbb{A}_{(j)}$ ,  $j = 1, 2, 3$ , in the calculation of  $I_2(\mathbf{b})$ , where  $\mathbf{b} = (4, 2)$ . Also, for particular values of  $\mathbf{a}|_{(0,2]} \in \mathbb{A}_{(1)}$ , namely,  $\mathbf{a}_2 = (2.5, 1)$  and  $\mathbf{a}(0, 2] = (5, 2)$  Figure 2(b) shows the workload paths  $\mathbf{W}^{(1)}$ . The same figure also displays examples corresponding to arrival paths in  $\mathbb{A}_{(j)}$ ,  $j = 2, 3$ . For the specific example from  $\mathbb{A}_{(1)}$  the figure shows that  $\mathbf{W}_1^{(1)} = \mathbf{a}_2 = (2.5, 1)$  and  $\mathbf{W}_0^{(1)} = \mathbf{a}(0, 2] - (1, 0) = (4, 2)$ .

This example underlines the difficulty in finding the rate function even for small timescales. We expect that the number of constrained sets like  $\mathbb{A}_{(j)}$  will grow exponentially with time duration  $t$ . However, the example gives us some insight on how to find some simple upper and lower bounds to  $I_t$  for any  $t \in \mathbb{N}$ .

### E. Bounds on $I_t$

In this subsection, we find simple expressions that give lower or upper bounds to  $\inf_{\mathbf{x} \in \mathbb{A}(u, \mathbf{b})} I_u^\sharp(\mathbf{x})$ , which in turn give the bounds to  $I_t$  and  $J$ . We focus on  $K = 2$  but similar results can be obtained for general  $K$ .

*Lemma 10:* For  $K = 2$ ,  $\mathbf{b} \in \mathbb{R}_+^2$ ,  $I_t(\mathbf{b})$  can be bounded as

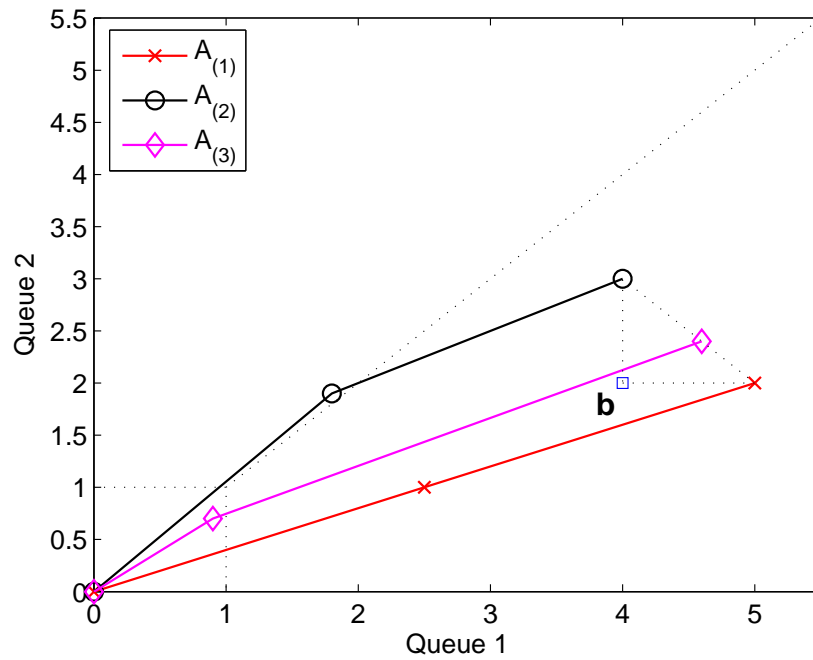
$$I_t(\mathbf{b}) \geq \min_{u \in (0, t]} u \sum_{k=1}^K \Lambda_1^* \left( \frac{1}{u} \left( \text{Proj}_{\mathbb{X}(u, \mathbf{b})}(\mathbf{0}) \right)^k \right) \quad (16)$$

and when  $\mathbf{b} \notin [0, 1]^2$ ,

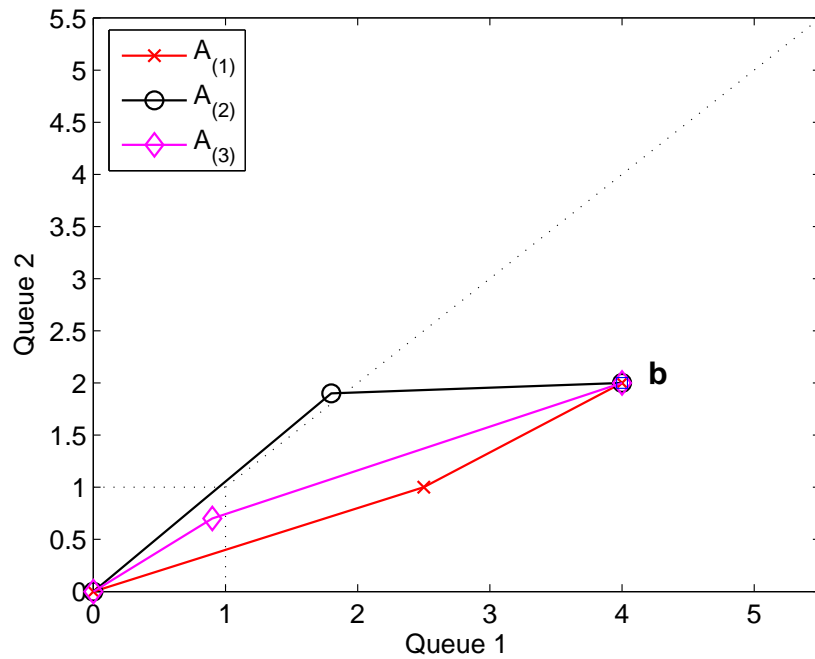
$$I_t(\mathbf{b}) \leq \min_{u \in (0, t]} u \sum_{k=1}^K \Lambda_1^* \left( \frac{1}{u} (b^k + (u-1)H(\mathbf{b})^k) \right), \quad (17)$$

where the convex set  $\mathbb{X}(u, \mathbf{b}) \subseteq \mathbb{R}_+^2$  is defined as

$$\mathbb{X}(u, \mathbf{b}) := \{\mathbf{b} + \mathbf{v} : v^1 + v^2 = (u-1), v^1, v^2 \geq 0\}. \quad (18)$$



(a) Trajectories of Arrival Paths



(b) Trajectories of Workload Paths

Fig. 2. Example of accumulated arrival and workload paths for calculation of  $I_2(\mathbf{b})$ .

*Proof:* See Appendix D. ■

Next we look at the tightness of the above bounds for an example of compound Poisson source process with exponential packet size. We expect the tightness to depend on the traffic load.

### F. Numerical Examples

Here we compare the finite-horizon rate functions  $I_2$  for three schedulers, namely, max-weight, GPS with equal weights, and a priority scheduler that gives higher priority to queue 1. We also examine the tightness of the bounds given in Lemma 10. Both of these are for an average of *i.i.d.* compound Poisson source processes with exponential packet size where the packet arrivals follow Poisson distribution of rate  $\lambda$  and the average packet size is  $1/\mu$  (see [40]). The function  $\Lambda^*$  for this process is given by

$$\Lambda^*(x) = \mu(\sqrt{x} - \sqrt{\lambda})^2,$$

for  $x \in \mathbb{R}_+$ . Note that this has domain  $\mathbb{R}_+$  and therefore satisfies Proposition 1. We once again make the simplifying assumption that the processes for the different users are statistically identical and that the rate-region is the unit simplex.

First we present results for the comparison of the different scheduling policies. Here we set  $\lambda = 0.1$  or  $0.3$  with the average packet size of  $1/\mu = 100$ . Fig. 3(a) and 3(b) show the finite-horizon (two-timestep) rate function  $I_2(\mathbf{b})$  for  $\lambda = 0.1$  and  $0.3$  respectively. However, these calculations are best appreciated when we compare them to the rate functions of other well-known scheduling policies. Fig. 4(a) compares the max-weight scheduler with a GPS scheduler with equal weights both at  $\lambda = 0.3$ . One can see that the rate function for the max-W scheduler is greater than the rate-function of the GPS scheduler for some range of  $\mathbf{b}$ , i.e., where  $b^1$  is much greater than  $b^2$  and *vice versa*. This means that for this range of  $\mathbf{b}$  and in the many-sources asymptotic sense, the work-conserving MW scheduler performs better than the work-conserving GPS scheduler (when we only consider the two time-step workload). The reason is that with  $b^1$  much greater than  $b^2$ , MW can serve queue 1 at full capacity (1), where GPS has to serve both queues equally at  $1/2$ . Hence, the workload under MW is less likely to reach  $(b^1, b^2)$  at the end of the two time-slots when the system starts from being empty. Similarly 4(b) compares the rate-functions for the max-weight scheduler and a priority scheduler where user 1 has higher priority. Since the max-weight policy does not discriminate between the two users, we find that

it is less likely for  $b^2$  to be large for the max-weight policy in comparison to the priority policy with the reverse being true for  $b^1$ .

Next we compare our bounds for the rate function of the max-weight policy with the exact rate function in the scenarios where we can calculate it by brute force. Figure 5(a) shows the upper and lower bounds and the actual values of  $I_t$ , for  $t = 10$ , at  $\mu = 0.01$ , and various values of  $\lambda = 0.1, 0.2, 0.3$  and when  $\mathbf{b} = (b^1, b^2 = 1)$  for various values of  $b^1$ . Figure 5(b) shows the corresponding minimizing  $t^*$  for the bounds and the actual expression of  $I_t$ . We note that for all  $\mathbf{b}$  in this example,  $J(\mathbf{b})$  is actually equal to  $I_t(\mathbf{b})$  for  $t = 10$  since all optimizing  $t^*$  is less than 10 (see (14)). This example shows that, in the range of  $\mathbf{b}$  in consideration, both bounds are tight and almost coincide when the traffic load is small, i.e.,  $\lambda = 0.1$ . However, when the traffic load is higher, the lowerbound becomes loose while the upperbound is still considerably tight.

It is interesting to note the optimal timescale  $t^*$  which the queues most likely to take to reach the level  $\mathbf{b}$ . Figure 5(b) shows that, for example, it is most likely to take only two timeslots for CPE process with  $\lambda = 0.2$  to reach the buffer level  $\mathbf{b} = (3, 1)$ , while the most likely timescale is four timeslots when the traffic load is higher ( $\lambda = 0.3$ ). Figure 6(c) and Figure 6(d) show the optimal trajectories of the accumulated arrival process and the workload process for  $\lambda = 0.2$  and 0.3, respectively.

Note that despite a potential exponential growth in computation of  $I_t$ , such computations commonly reduce to simpler cases. For instance, consider the calculation of  $I_{10}$  in Figs. 6(c) and 6(d), in which we sequentially calculate  $I_1, I_2, I_3$ . However, this sequential computation stops as soon as one reaches the optimal timescale  $t^* < t$ . For the values of the vector  $b$  in Figs. 6(c) and 6(d), for instance,  $t^*$  was at most 4, making the calculations of  $I_6-I_{10}$  unnecessary.

### VIII. CONCLUDING REMARKS

In this paper, we have established a many-sources LDP for the stationary (infinite-horizon) workload for multi-queue single-server system with simplex rate region and under maximum-weight scheduling, when the arrival processes assumed to satisfy certain many-sources sample path LDP. To extend the LDP of the arrival processes to the LDP of the workloads, we employed Garcia's extended contraction principle, which applies to quasi-continuous mappings. Along the way, we also establish an LDP for the finite-horizon workload in a very general setting of arbitrary compact, convex, and coordinate convex rate region under max-weight scheduling. We

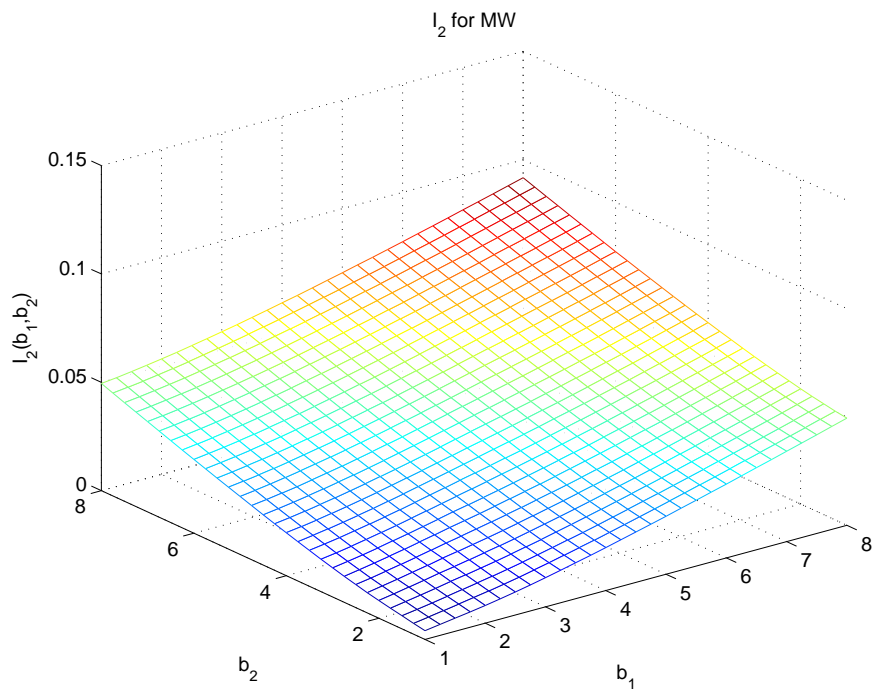
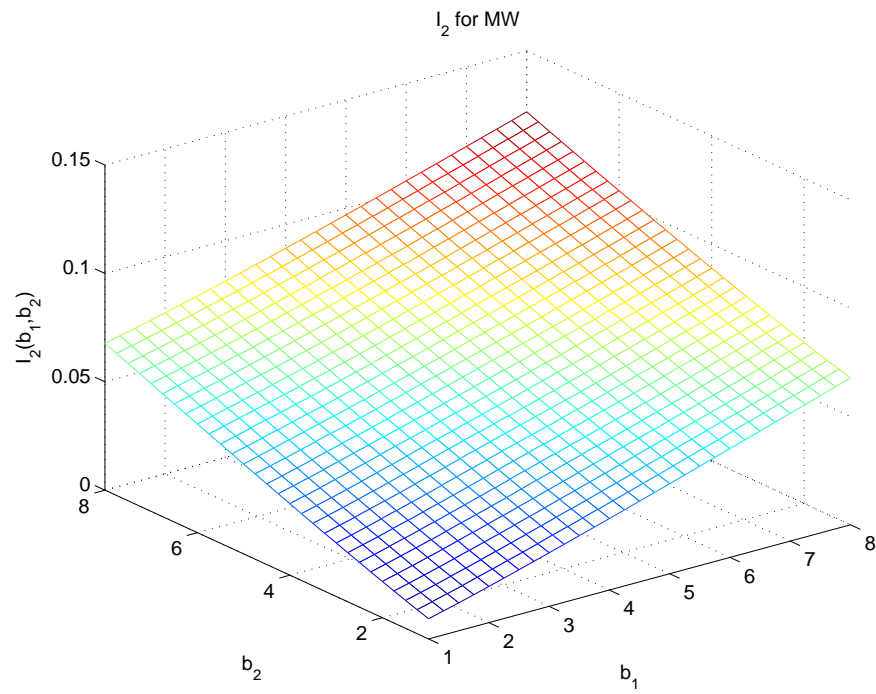


Fig. 3. Finite-horizon rate function  $I_2(\mathbf{b})$  for the max-weight scheduler.

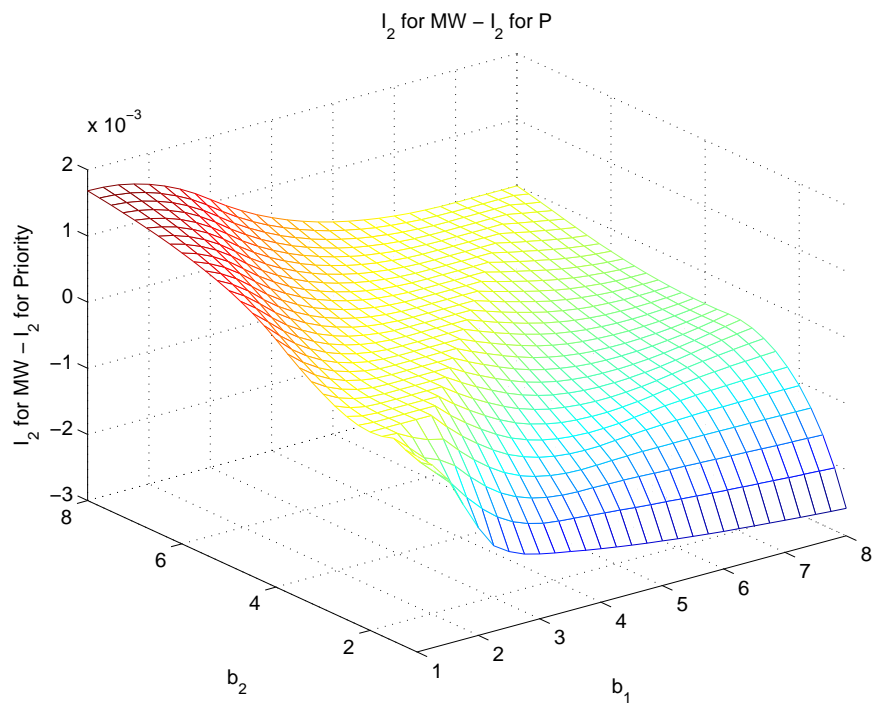
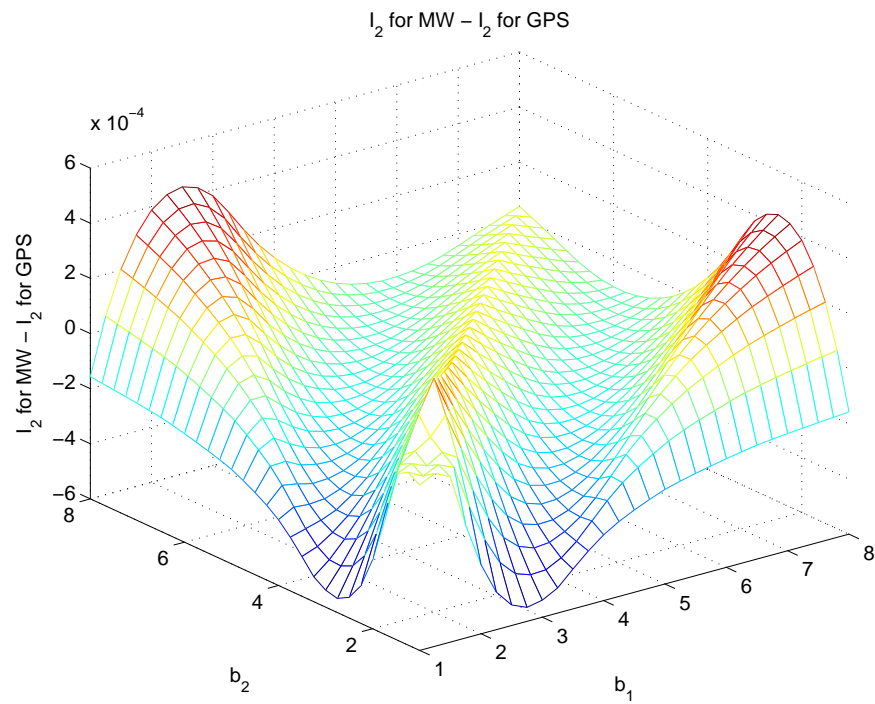
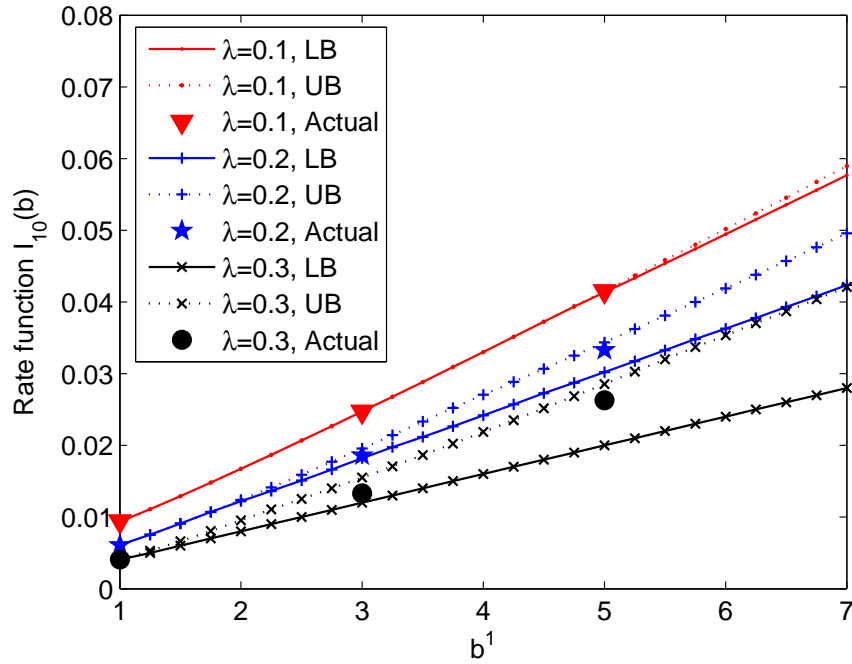
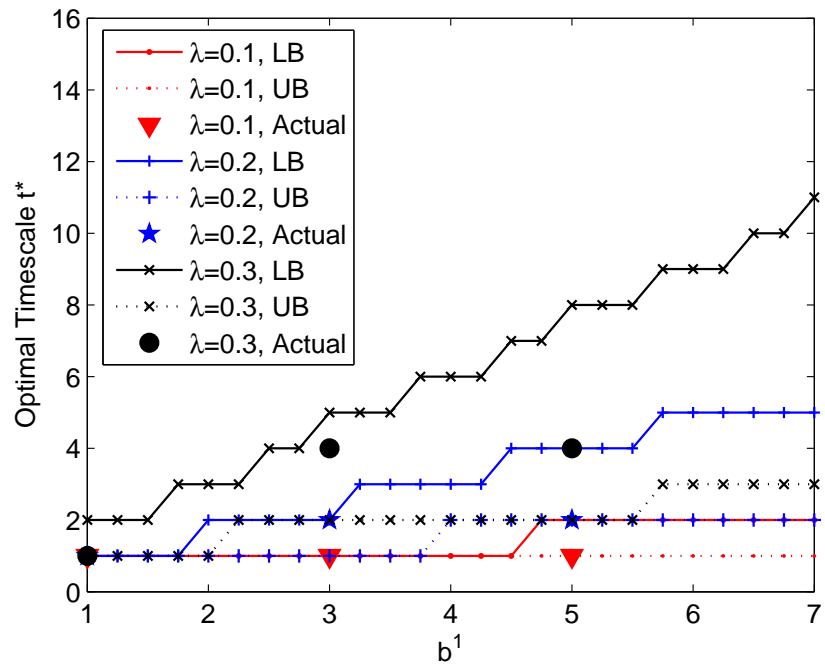


Fig. 4. Comparison of the rate-functions of the max-weight scheduler, the GPS scheduler with equal weights and the priority scheduler where user 1 has higher priority.

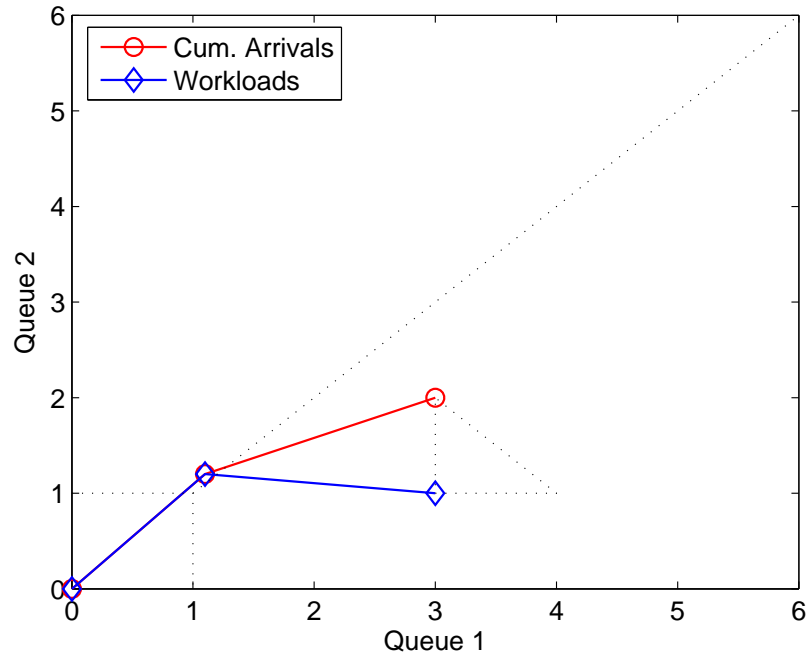
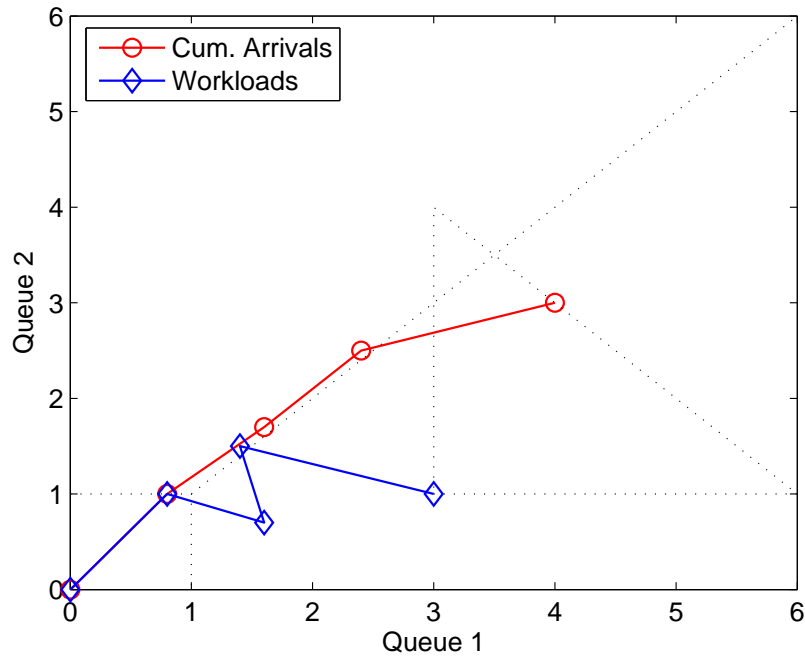


(a) Rate functions



(b) Optimal Timescales

Fig. 5. Example of the rate function  $I_{10}(\mathbf{b})$  and its upper and lower bounds, and their corresponding optimizing  $t^*$  and optimal trajectories, when  $\mathbf{b} = (b^1, b^2 = 1)$ .

(c) Optimal Trajectory for  $\mathbf{b} = (3, 1)$  and  $\lambda = 0.2$ (d) Optimal Trajectory for  $\mathbf{b} = (3, 1)$  and  $\lambda = 0.3$ Fig. 6. The optimal accumulated arrival and workload trajectories when  $\mathbf{b} = (3, 1)$  at  $\lambda = 0.2$  and  $\lambda = 0.3$ , respectively.

gave the associated rate functions and the expression of the infinite-horizon rate function in term of the finite-horizon ones, when the arrival processes have *i.i.d.* increments.

Next, we catalogue some interesting areas of future research. The extension of our LDP result for the infinite-horizon workload in the case of an arbitrary compact, convex, and coordinate convex rate region remains open by and large. The main difficulty in establishing an LDP for the infinite-horizon workload is in showing the quasi-continuity of the infinite-horizon workload mapping. In the case of simplex rate region, the infinite-horizon workload mapping was shown to be reducible to a finite-horizon mapping whose quasi-continuity was established via induction. Another question that has only been partially explored in the current paper is the nature of Assumptions 2-3. In particular, beyond Proposition 1, the relationship between stochastic mixing properties of the arrival process and the analytical properties of  $I_t^\#$  and  $I^\#$  remains an important area of future research.

## APPENDIX A

### PROOF OF LEMMA 1

Lemma 1 states that it is possible to construct quasi-continuous functions  $H$  and  $H^{\text{wc}}$  such that

$$H(\mathbf{W}_t) \in \arg \max_{\mathbf{R} \in \mathcal{R}} \langle \mathbf{R}, \mathbf{W}_t \rangle,$$

and

$$H^{\text{wc}}(\mathbf{W}_t) = \begin{cases} \text{Proj}_{\mathcal{R}}(\mathbf{W}_t) & \mathbf{W}_t \in \prod_{k=1}^K [0, C_k) \\ H(\mathbf{W}_t) & \text{otherwise} \end{cases}.$$

We do this using selection theorems (see [58] and [59]) for correspondences associated with max-weight scheduling:

$$\mathcal{H}(\mathbf{W}_t) := \arg \max_{\mathbf{R} \in \mathcal{R}} \langle \mathbf{R}, \mathbf{W}_t \rangle. \quad (19)$$

First we define the following analytical properties of correspondences.

*Definition 5 (Local-boundedness):* A correspondence  $F : \mathcal{X} \rightrightarrows \mathcal{Y}$  is locally-bounded (*lb*) [56, Defn. 5.15, pp. 157–158] at point  $x \in \mathcal{X}$  if there exists a neighbourhood  $V$  of  $x$  such that  $F(V) := \cup_{a \in V} F(a)$  is bounded in  $\mathcal{Y}$ . Correspondence  $F$  is deemed *lb* if it is *lb* for every point  $x \in \mathcal{X}$ .

*Definition 6 (Outer semicontinuity):* For a correspondence  $F : \mathcal{X} \rightrightarrows \mathcal{Y}$  define the outer-limit at  $x \in \mathcal{X}$  to be

$$\limsup_{a \rightarrow x} F(a) := \cup_{x^n \rightarrow x} \limsup_{n \rightarrow \infty} F(x^n) = \{y | \exists x^n \rightarrow x, \exists y^n \rightarrow y \text{ with } y^n \in F(x^n)\}.$$

Then  $F$  is outer semicontinuous (*osc*) [56, Defn. 5.4, pg. 152] at  $x \in \mathcal{X}$  if  $\limsup_{a \rightarrow x} F(a) \subset F(x)$ . Correspondence  $F$  is deemed *osc* if it is *osc* for every point  $x \in \mathcal{X}$ .

*Definition 7 (Upper semicontinuity):* A correspondence  $F : \mathcal{X} \rightrightarrows \mathcal{Y}$  is upper semicontinuous (also known as upper hemicontinuous, [56, Thm. 5.19, pg. 160] and [57]) at point  $x \in \mathcal{X}$  if for any open neighbourhood  $V$  of  $F(x)$ , there exists a neighbourhood  $U$  of  $x$  such that  $F(a) \subseteq V$  for all  $a \in U$ . Alternatively,  $F$  is said to be upper semicontinuous at  $x \in \mathcal{X}$  if whenever we have sequences  $\{x_m\} \in \mathcal{X}$  and  $\{y_m\} \in \mathcal{Y}$  such that  $y_m \in F(x_m)$  for all  $m$ ,  $x_m \rightarrow x$  and  $y_m \rightarrow y$ , then  $y \in F(x)$ . Correspondence  $F$  is deemed *usc* if it is *usc* for every point  $x \in \mathcal{X}$ .

Now we can prove Lemma 1.

*Proof of Lemma 1:* The scheduling function  $\mathcal{H}(\cdot) : \mathbb{R}_+^K \mapsto \mathcal{P}(\mathcal{R})$  is a maximal monotone correspondence [56, Thm. 12.17, pp. 542–543] that picks closed and convex subsets of  $\mathcal{R}$  for every  $\mathbf{x} \in \mathbb{R}_+^K$ ; thus, compact subsets of  $\mathbb{R}_+^K$ . It is, therefore, both *lb* and *usc* from [56, Ex. 12.8b, pg. 536]. Now it follows that we get a quasi-continuous selection by [58, Thm. 2] and [59, Thms. 2.4 & 2.5] since  $H(\mathbf{W}_t(\cdot))$  is *usc* and has compact values.

To prove that  $H^{\text{wc}}$  is a quasi-continuous selection we first use the same steps as above but with a restricted domain, i.e., for  $\mathcal{H}(\cdot) : \mathbb{R}_+^K \setminus \prod_{k \in \mathcal{K}} [0, C_k) \mapsto \mathcal{P}(\mathcal{R})$  since  $\mathbb{R}_+^K \setminus \prod_{k \in \mathcal{K}} [0, C_k)$  is a Baire space [42] so that results from [58], [59] still apply. From the definition of the work-conserving max-weight scheduler, if  $\mathbf{x} \in \prod_{k \in \mathcal{K}} [0, C_k)$ , then we get continuity from within this set from the properties of  $\text{Proj}_{\mathcal{R}}(\cdot)$ . Thus, we can satisfy the definition of quasi-continuity from Defn. 1. ■

We refer the reader to [60, Thm. 2.2] and [61, Thm. 3.4] for an exposition and for generalizations of the result from [58], [59].

## APPENDIX B

### PROOF OF LEMMAS 2-3

Next we prove Lemmas 2 which uses the following fact:

*Fact 2:* Assume  $\mathcal{X}, \mathcal{Y}$  are metric spaces and  $F_1$  and  $F_2$  are functions from  $\mathcal{X}$  onto  $\mathcal{Y}$ . If  $F_1$  is quasi-continuous at  $x \in \mathcal{X}$  and  $F_2$  is continuous at  $x$ , then  $F_1 + F_2$  is quasi-continuous at  $x$ .

Below we state and prove Lemma 2 for a simplex rate region  $\mathcal{R}_s$  to avoid unnecessary notational complexities; the result readily extends to a general rate region as discussed in Remark 9.

*Lemma 2:* For  $t \in \mathbb{N}$ ,  $G_t^{\text{wc}}$  for  $\mathcal{R}_s$  is quasi-continuous on  $\mathbb{R}_+^{Kt}$  with respect to the scaled uniform norm topology.

*Proof:* Using our queueing equation we first observe the following recursive relation between  $G_t^{\text{wc}}$  and  $G_{t-1}^{\text{wc}}$  for any  $t \in \{2, 3, \dots\}$  and  $\mathbf{x} \in \mathcal{D}_\mu^K$ :

$$G_t^{\text{wc}}(\mathbf{x}|_{(0,t]}) = \mathbf{x}_1 + [G_{t-1}^{\text{wc}}(\mathbf{x}|_{(1,t]}) - H^{\text{wc}}(G_{t-1}^{\text{wc}}(\mathbf{x}|_{(1,t]}))]^+, \quad (20)$$

where we used the fact that  $\mathbf{W}_0^{\text{wc}}(\mathbf{x}|_{(0,t]}) = G_t^{\text{wc}}(\mathbf{x}|_{(0,t]})$ , and  $\mathbf{W}_1^{\text{wc}}(\mathbf{x}|_{(1,t]}) = G_{t-1}^{\text{wc}}(\mathbf{x}|_{(1,t]})$  when the initial backlog at time  $-t$  is  $\mathbf{0}$ .

Equation (20) says that  $G_t^{\text{wc}}(\mathbf{x}|_{(0,t]})$  depends linearly on  $\mathbf{x}_1$ . This implies the following simple but consequential observations:

*Observation 1:* If  $G_t^{\text{wc}}$  is quasi-continuous at  $\mathbf{x}|_{(0,t]}$ , then it is quasi-continuous at  $\tilde{\mathbf{x}}|_{(0,t]} := (\tilde{\mathbf{x}}_1, \mathbf{x}_2, \dots, \mathbf{x}_t)$  for any  $\tilde{\mathbf{x}}_1 \in \mathbb{R}_+^K$ , and if  $G_t^{\text{wc}}$  is continuous at  $\mathbf{x}|_{(0,t]}$ , then it is also continuous at  $\tilde{\mathbf{x}}|_{(0,t]}$ .

*Observation 2:* If  $G_t^{\text{wc}, \mathcal{R}_s}(\mathbf{x}^n|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{x}|_{(0,t]})$  for a sequence  $\{\mathbf{x}^n|_{(0,t]}\}$  such that  $\mathbf{x}^n|_{(0,t]} \rightarrow \mathbf{x}|_{(0,t]}$ , then for any sequence  $\{\tilde{\mathbf{x}}^n|_{(0,t]} = (\tilde{\mathbf{x}}_1^n, \mathbf{x}_2^n, \dots, \mathbf{x}_t^n)\}$  where  $\tilde{\mathbf{x}}_1^n \rightarrow \mathbf{x}_1$ , we also have  $G_t^{\text{wc}}(\tilde{\mathbf{x}}^n|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{x}|_{(0,t]})$ .

Using the recursive relation in (20), we prove this lemma by induction on  $t \in \mathbb{N}$ . For every  $t \in \mathbb{N}$  we assume that the system is empty at  $-t$ . Therefore, the arrival sample-path prior to  $-t$  has no influence on  $\mathbf{W}_{0,t}^{\text{wc}}$ . Thus, we will only specify the values of the constructed sequences up until time  $t$ ; the extension to sample-paths in  $\mathcal{D}_\mu^K$  while ensuring that the system is empty at  $-t$  is trivial. For  $t = 1$ ,  $G_1^{\text{wc}}(\mathbf{a}_1) = \mathbf{a}_1$ , hence  $G_1^{\text{wc}}$  is continuous on  $\mathbb{R}_+^K$ . Assuming that  $G_t^{\text{wc}}$  is quasi-continuous on  $\mathbb{R}_+^{Kt}$ , we want to show that  $G_{t+1}^{\text{wc}}$  is quasi-continuous on  $\mathbb{R}_+^{K(t+1)}$ . Using the fact that the  $[\cdot]^+$  function is continuous, Remark 2, and Fact 2, it suffices to show that the function  $F_t := G_t^{\text{wc}} - H^{\text{wc}} \circ G_t^{\text{wc}}$ , is quasi-continuous on  $\mathbb{R}_+^{Kt}$  to show that  $G_{t+1}^{\text{wc}}$  is quasi-continuous. In particular, for any arrival sample path  $\mathbf{a} \in \mathcal{D}_\mu^K$ , we need to show that  $F_t$  is quasi-continuous at  $\mathbf{a}|_{(0,t]}$ . It suffices to show that it is possible to select a sequence  $\hat{\mathbf{a}}^n \rightarrow \mathbf{a}$  (in  $\mathcal{D}_\mu^K$ ) for which

$$G_t^{\text{wc}}(\hat{\mathbf{a}}^n|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{a}|_{(0,t]}), \quad (21)$$

$$H^{\text{wc}} \circ G_t^{\text{wc}}(\hat{\mathbf{a}}^n|_{(0,t]}) \rightarrow H^{\text{wc}} \circ G_t^{\text{wc}}(\mathbf{a}|_{(0,t]}), \quad (22)$$

such that both  $G_t^{\text{wc}}(\cdot)$  and  $H^{\text{wc}} \circ G_t^{\text{wc}}(\cdot)$  are continuous at every  $\hat{\mathbf{a}}^n|_{(0,t]}$ . Note that in contrast to Fact 2 we are adding two quasi-continuous and showing that the sum is quasi-continuous; the key to our proof is to ensure that we use the same sequence for both functions.

We show this by first noting that the induction hypothesis, i.e., quasi-continuity of  $G_t^{\text{wc}}$ , and the definition of quasi-continuity ensure that there exists a sequence  $\{\mathbf{a}^n\}$  such that  $\mathbf{a}^n \rightarrow \mathbf{a}$ , in the scaled uniform norm topology, such that  $G_t^{\text{wc}}(\mathbf{a}^n|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$ , and  $G_t^{\text{wc}}(\cdot)$  is continuous at  $\mathbf{a}^n|_{(0,t]}$  for all  $n$ . We will construct the desired sequence  $\{\hat{\mathbf{a}}^n\}$  by modifying  $\hat{\mathbf{a}}_1^n$  appropriately. We proceed by considering the following two cases, depending on the value of  $\mathbf{a}_1$ .

**Case 1:**  $\mathbf{a}_1 > \mathbf{0}$ , i.e., every component of the  $\mathbf{a}_1 \in \mathbb{R}_+^K$  is positive. Let  $\epsilon > 0$  be the smallest component of  $\mathbf{a}_1$ , i.e.,  $\epsilon = \min_{k \in \mathcal{K}} a_1^k$ . Since  $H(\cdot)$  is quasi-continuous, it is possible to choose a sequence of (workload) vectors  $\{\mathbf{w}^n\}$  such that  $\mathbf{w}^n \rightarrow G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$ , and  $H$  is continuous at  $\mathbf{w}^n$  for all  $n$ . Now, we define

$$\begin{aligned} \tilde{\mathbf{a}}_1^n &:= \mathbf{w}^n - [F_{t-1}(\mathbf{a}^n|_{(1,t]})]^+ = \mathbf{w}^n - G_t^{\text{wc}}(\mathbf{a}^n|_{(0,t]}) + \mathbf{a}_1^n \\ &= (\mathbf{w}^n - G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})) + (G_t^{\text{wc}}(\mathbf{a}|_{(0,t]}) - G_t^{\text{wc}}(\mathbf{a}^n|_{(0,t]})) + (\mathbf{a}_1^n - \mathbf{a}_1) + \mathbf{a}_1. \end{aligned} \quad (23)$$

It is clear from the last relationship that  $\tilde{\mathbf{a}}_1^n \rightarrow \mathbf{a}_1$ . We still need to ensure that  $\tilde{\mathbf{a}}_1^n \geq \mathbf{0}$  since negative quantities are involved in the definition. We do this by using the facts that every component of  $\mathbf{a}_1 \in \mathbb{R}^K$  is greater than or equal to  $\epsilon > 0$ , and that  $\mathbf{w}^n \rightarrow G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$ ,  $G_t^{\text{wc}}(\mathbf{a}^n|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$ , and  $\mathbf{a}_1^n \rightarrow \mathbf{a}_1$ . These facts imply that there exists an  $n_\epsilon$  such that for all  $n > n_\epsilon$  we have  $\|\mathbf{w}^n - G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})\| < \epsilon/3$ ,  $\|G_t^{\text{wc}, \mathcal{R}_s}(\mathbf{a}|_{(0,t]}) - G_t^{\text{wc}}(\mathbf{a}^n|_{(0,t]})\| < \epsilon/3$  and  $\|\mathbf{a}_1^n - \mathbf{a}_1\| < \epsilon/3$  (with the square norm) which then together with (23) imply that, for the sequence  $\tilde{\mathbf{a}}_1^{m+n_\epsilon}$ , we always have non-negativity of all components. Hence, we construct a new sequence  $\{\hat{\mathbf{a}}_{(0,t]}^n\}$  where  $\hat{\mathbf{a}}_1^n = \tilde{\mathbf{a}}_1^{n+n_\epsilon}$  and  $\hat{\mathbf{a}}_{(1,t]}^n = \mathbf{a}_{(1,t]}^{n+n_\epsilon}$ .

This new sequence  $\hat{\mathbf{a}}_{(0,t]}^n$  is the sequence we are after because using the induction hypothesis together with Observations 1 and 2, we have that  $G_t^{\text{wc}}(\hat{\mathbf{a}}^n|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$ , and  $G_t^{\text{wc}}$  is continuous at  $\hat{\mathbf{a}}^n|_{(0,t]}$  for all  $n$ . Furthermore, by construction

$$G_t^{\text{wc}}(\hat{\mathbf{a}}^n|_{(0,t]}) = \hat{\mathbf{a}}_1^n + [F_{t-1}(\hat{\mathbf{a}}^n|_{(1,t]})]^+ = \tilde{\mathbf{a}}_1^{n+n_\epsilon} + [F_{t-1}(\mathbf{a}^{n+n_\epsilon}|_{(1,t]})]^+ = \mathbf{w}^{n+n_\epsilon}. \quad (24)$$

Hence, we have shown that there exists a sequence  $\hat{\mathbf{a}}_{(0,t]}^n$  satisfying (21) and (22). In addition, the continuity of  $H^{\text{wc}} \circ G_t^{\text{wc}}$  at  $\hat{\mathbf{a}}^n|_{(0,t]}$  for all  $n$  is a direct consequence of continuity of  $G_t^{\text{wc}}$  at  $\hat{\mathbf{a}}^n|_{(0,t]}$  and continuity of  $H^{\text{wc}}$  at  $\mathbf{w}^{n+n_\epsilon}$ , which is equal to  $G_t^{\text{wc}}(\hat{\mathbf{a}}^n|_{(0,t]})$ , for all  $n$ .

**Case 2:**  $\mathbf{a}_1 \geq \mathbf{0}$ . Let  $\mathcal{K}_1 := \{k \in \mathcal{K} : a_1^k = 0\}$  and let  $\mathcal{K}_2 := \arg \max_{k \in \mathcal{K}} \mathbf{W}_0^k(\mathbf{a}|_{(0,t)})C_k$ . Without loss of generality, by permuting the user labels we can assume that the first  $\hat{K} := |\mathcal{K}_1 \cup \mathcal{K}_2|$  components of  $\mathbf{a}_1$  are either in  $\mathcal{K}_1$  or in  $\mathcal{K}_2$ ; thus, the rest of the  $K - \hat{K}$  components are both positive and not part of the scheduling decision made at time 0 with arrival sequence  $\mathbf{a}_{(0,t]}$ . Now consider the sequence  $\mathbf{a}_1^m := [1/mC]_{\hat{K}} + \mathbf{a}_1$  where  $[1/mC]_{\hat{K}}$  is short-hand for a vector with  $1/(mC_k)$  in the first  $\hat{K}$  components and 0 in the remaining coefficients; by construction  $\mathbf{a}_1^m$  converges to  $\mathbf{a}_1$  such that for every  $m$  every component of  $\mathbf{a}_1^m$  is positive. We construct a sequence  $\{\mathbf{a}^m|_{(0,t]}\}$  with this  $\mathbf{a}_1^m$  and  $\mathbf{a}^m|_{(1,t]} = \mathbf{a}|_{(1,t]}$ . It is obvious that  $G_t^{\text{wc}}(\mathbf{a}^m|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$  since

$$G_t^{\text{wc}}(\mathbf{a}^m|_{(0,t]}) = \mathbf{a}_1^m + [F_{t-1}(\mathbf{a}^m|_{(1,t]})]^+ = \mathbf{a}_1^m + [F_{t-1}(\mathbf{a}|_{(1,t]})]^+ = [1/mC]_{\hat{K}} + G_t^{\text{wc}}(\mathbf{a}|_{(0,t]}).$$

When  $G_t^{\text{wc}, \mathcal{R}_s}(\mathbf{a}|_{(0,t]}) \notin [0, C)^K$ , by construction, we have

$$H^{\text{wc}} \circ G_t^{\text{wc}}(\mathbf{a}^m|_{(0,t]}) = H^{\text{wc}}(G_t^{\text{wc}}(\mathbf{a}^m|_{(0,t]})) = H^{\text{wc}}(G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})) = H \circ G_t^{\text{wc}}(\mathbf{a}|_{(0,t]}),$$

where the function  $H(\cdot)$  is the regular max-weight scheduling function (with  $\mathcal{R}_s$ ). On the other hand, if  $G_t^{\text{wc}}(\mathbf{a}|_{(0,t]}) \in [0, C)^K$ , then the continuity of  $\text{Proj}_{\mathcal{R}_s}(\cdot)$  yields  $H^{\text{wc}} \circ G_t^{\text{wc}}(\mathbf{a}^m|_{(0,t]}) \rightarrow H^{\text{wc}} \circ G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$ .

Since for each  $m$  we have that  $\mathbf{a}_1^m$  has all elements strictly positive, we can use the construction from **Case 1** but with  $\mathbf{a}^m|_{(0,t]}$  in place of  $\mathbf{a}|_{(0,t]}$ . In particular, for each  $m$ , we can now generate a sequence  $\{\tilde{\mathbf{a}}^{m,n}\}$  such that  $\tilde{\mathbf{a}}_1^{m,n} \rightarrow \mathbf{a}_1^m$  as  $n \rightarrow +\infty$ ,  $\tilde{\mathbf{a}}^{m,n}|_{(1,t]} = \mathbf{a}^n|_{(1,t]}$ , and by using Observations 1 and 2, the following hold

$$G_t^{\text{wc}}(\tilde{\mathbf{a}}^{m,n}|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{a}^m|_{(0,t]}), \quad (25)$$

$$H^{\text{wc}} \circ G_t^{\text{wc}}(\tilde{\mathbf{a}}^{m,n}|_{(0,t]}) \rightarrow H^{\text{wc}} \circ G_t^{\text{wc}}(\mathbf{a}^m|_{(0,t]}), \quad (26)$$

with both  $G_t^{\text{wc}}(\cdot)$  and  $H^{\text{wc}} \circ G_t^{\text{wc}}(\cdot)$  being continuous at  $\tilde{\mathbf{a}}^{m,n}|_{(0,t]}$  for all  $m, n$ .

Now we define the sequence  $\hat{\mathbf{a}}^m = \tilde{\mathbf{a}}^{m,m}$  as the sequence we are after. By construction, we have  $\hat{\mathbf{a}}^m|_{(0,t]} \rightarrow \mathbf{a}|_{(0,t]}$  and both  $G_t^{\text{wc}}(\cdot)$  and  $H^{\text{wc}} \circ G_t^{\text{wc}}(\cdot)$  continuous at all  $\hat{\mathbf{a}}^m|_{(0,t]}$ . Since  $G_t^{\text{wc}}(\mathbf{a}^m|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$  and  $H^{\text{wc}} \circ G_t^{\text{wc}}(\mathbf{a}^m|_{(0,t]}) \rightarrow H^{\text{wc}} \circ G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$ , it follows from (25) and (26) that  $G_t^{\text{wc}}(\hat{\mathbf{a}}^m|_{(0,t]}) \rightarrow G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$  and  $H^{\text{wc}, \mathcal{R}_s} \circ G_t^{\text{wc}}(\hat{\mathbf{a}}^m|_{(0,t]}) \rightarrow H^{\text{wc}} \circ G_t^{\text{wc}}(\mathbf{a}|_{(0,t]})$ . ■

*Remark 9:* The proof above can be carried out for every rate-region in the class that we are interested in. The argument presented in **Case 1** would remain exactly the same but the argument presented in **Case 2** would have to be modified to account for a further characterization of

$\arg \max_{\mathbf{R} \in \mathcal{R}} \langle \mathbf{W}_0(\mathbf{a}|_{(0,t]}), \mathbf{R} \rangle$ , especially when it is not a singleton. The components of  $\mathbf{a}_1$  will need to be adjusted in such a manner so as to not perturb  $\arg \max_{\mathbf{R} \in \mathcal{R}} \langle \mathbf{W}_0(\mathbf{a}^m|_{(0,t]}), \mathbf{R} \rangle$  for the adjusted sequence  $\mathbf{a}^m|_{(0,t]}$ . As the case of a non-singleton  $\arg \max_{\mathbf{R} \in \mathcal{R}} \langle \mathbf{W}_0(\mathbf{a}^m|_{(0,t]}), \mathbf{R} \rangle$  will correspond to a specific set of values of  $\mathbf{W}_0(\mathbf{a}|_{(0,t]})$  (a cone) such that the boundary of the rate-region and a hyper-plane intersect at more than one point, we will need to use the normal corresponding to the hyper-plane in constructing the appropriate perturbation. Thus, on a case-by-case basis the same argument can be carried out for every rate-region.

Finally, we prove Lemma 3.

*Lemma 3:* For  $t \in \mathbb{N}$ ,  $G_t(\cdot)$  is quasi-continuous on  $\mathbb{R}_+^{K \times t}$ .

*Proof:* Consider the queueing equation, i.e.,

$$\mathbf{W}_{T-1} = [\mathbf{W}_T - \mathbf{R}_T]^+ + \mathbf{a}_T \quad T \in \mathbb{N},$$

where  $\mathbf{R}_{(\cdot)} \in \mathcal{H}(\mathbf{W}_{(\cdot)})$  where  $\mathcal{H}(\cdot)$  is the maximal monotone correspondence defined in (19). Assume that we start the system at (fixed) time  $t \in \mathbb{N}$  with workload vector  $\mathbf{W}_t \in \mathbb{R}_+^K$ ; we will often assume that  $\mathbf{W}_t = \mathbf{0}$ . First, for  $t \in \mathbb{N}$  we define a correspondence  $\mathcal{G}_t(\mathbf{W}_t, \mathbf{a}|_{(0,t]}) : \mathbb{R}_+^{(t+1)K} \rightrightarrows \mathbb{R}_+^K$  that represents all possible workload vectors at time 0 that can be achieved from the inputs  $(\mathbf{W}_t, \mathbf{a}|_{(0,t]})$ . This results from the successive application of the queueing equation where we use all possible values in  $\mathcal{H}(\cdot)$  based upon the workload vector that results at each step. Our goal is show that  $\mathcal{G}_t(\cdot)$  admits a quasicontinuous selection which we shall call  $\tilde{G}_t(\cdot)$ . For  $t > 1$  it suffices to establish that

$$\begin{aligned} & \mathcal{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]}) - \mathcal{H}(\mathcal{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]})) := \\ & \{ \mathbf{x} \in \mathbb{R}^K : \mathbf{x} = \mathbf{y} - \mathbf{z} \text{ for some } \mathbf{y} \in \mathcal{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]}) \text{ and } \mathbf{z} \in \mathcal{H}(\mathcal{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]})) \} \end{aligned}$$

admits a quasicontinuous selection which we shall call (with an abuse of notation)  $\tilde{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]}) - H(\tilde{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]}))$ : 1) since  $[\cdot]^+$  is a continuous function, we have  $[\tilde{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]}) - H(\tilde{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]}))]^+$  being a quasi-continuous function; and finally, 2) by properties of projections and by the definition of quasi-continuity we get the quasi-continuity of  $\tilde{G}_t(\mathbf{W}_t, \mathbf{a}|_{(0,t]})$ . We should add a note of caution here that even though for notational convenience we write

$$\tilde{G}_t(\mathbf{W}_t, \mathbf{a}|_{(0,t]}) = [\tilde{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]}) - H(\tilde{G}_{t-1}(\mathbf{W}_t, \mathbf{a}|_{(1,t]}))]^+ + \mathbf{a}_1,$$

it need not be the case that the quasi-continuous selection that we obtain for  $\tilde{G}_t$  be related using the above queueing equation to the quasi-continuous selection for  $\tilde{G}_{t-1}$ . Additionally, we may not even use the quasi-continuous selection  $H(\cdot)$ . Therefore, the property highlighted in Remark 8 need not hold.

The proof will once again use mathematical induction where our induction step will assume that  $\mathcal{G}_{t-1}(\cdot)$  is *osc*. Since  $\mathbf{W}_0 \leq \mathbf{W}_t + \mathbf{a}(0, t]$  for any  $\mathbf{W}_0 \in \mathcal{G}_t(\mathbf{W}_t, \mathbf{a}|_{(0,t]})$ , it follows that  $\mathcal{G}_t(\mathbf{W}_t, \mathbf{a}|_{(0,t]})$  is *lb*. Now using [56, Ex. 12.8b, pg. 536] we know that  $\mathcal{H}(\cdot)$  is both *lb* and *osc*, and therefore by [56, Thm. 5.19, pg. 160] it is also *usc*. Then using [56, Prop. 5.52b, pp. 184–185] we have  $\mathcal{H}(\mathcal{G}_{t-1}(\cdot))$  being *osc*. Therefore it also follows that  $\mathcal{G}_{t-1}(\cdot) - \mathcal{H}(\mathcal{G}_{t-1}(\cdot))$  is also *osc*. Once this has been demonstrated the induction step is very easy as  $\mathcal{G}_t(\cdot)$  is obtained from  $\mathcal{G}_{t-1}(\cdot) - \mathcal{H}(\mathcal{G}_{t-1}(\cdot))$  by continuous transformations, as mentioned above. This same method also allows us to establish the initial step of the induction procedure; note that we will be dealing with  $\mathcal{H}(\mathbf{W}_t)$  and  $\mathbf{W}_t - \mathcal{H}(\mathbf{W}_t)$  in this case.

Since  $\mathcal{G}_{t-1}(\cdot) - \mathcal{H}(\mathcal{G}_{t-1}(\cdot))$  is *lb*, using [56, Ex. 12.8b, pg. 536] we have  $\mathcal{G}_{t-1}(\cdot) - \mathcal{H}(\mathcal{G}_{t-1}(\cdot))$  also being *usc*. Finally, we get a quasi-continuous selection by [60, Thm. 2.2] and [61, Thm. 3.4] since  $\mathcal{G}_{t-1}(\cdot) - \mathcal{H}(\mathcal{G}_{t-1}(\cdot))$  is *usc* and takes compact values.

The required result then follows by setting  $G_t(\mathbf{a}|_{(0,t]}) = \tilde{G}_t(\mathbf{0}, \mathbf{a}|_{(0,t]})$ . ■

## APPENDIX C

### PROOF OF LEMMAS 5-7

First we prove Lemma 5:

*Lemma 5:* Consider an arrival process  $\mathbf{a} \in \mathcal{D}_\mu^K$ . There exists a  $s^* = s^*(\mathbf{a}) < \infty$  such that the workloads at time  $-s^*$  under  $\mathbf{a}$  falls within the rate region  $\mathcal{R}_s$ , i.e.,  $G^{\text{wc}}(\mathbf{a}|_{(s^*, \infty)}) \in \mathcal{R}_s$ . Furthermore, for any sequence of arrival processes  $\{\mathbf{a}^n \in \mathcal{D}_\mu^K\}$  converging to  $\mathbf{a}$  (in scaled uniform topology), the workloads at time  $-s^*$  under  $\mathbf{a}^n$ , when  $n$  is large enough, also fall within the rate region  $\mathcal{R}_s$ , i.e.,  $\exists n_0$  such that  $G^{\text{wc}}(\mathbf{a}^n|_{(s^*, \infty)}) \in \mathcal{R}_s$  for  $n > n_0$ .

*Proof:* Consider the normalized sum arrivals and the normalized sum workloads, and follow the proof in [24], [26] for the (aggregate) single-queue scenario. Given the definition of  $H^{\text{wc}}$  and the simplex capacity region  $\mathcal{R}_s$ , the queue dynamics for the normalized sum workload is that of a single queue whose arrivals are the normalized sum of the arrivals, i.e.,

$$\hat{\mathbf{W}}_{t-1} = [\hat{\mathbf{W}}_t - 1]^+ + \hat{\mathbf{a}}_t, \quad (27)$$

where recall that the hat ( $\hat{\cdot}$ ) notation means the normalized sum over all users, i.e.  $\hat{\mathbf{a}}_t = \sum_{k=1}^K a_t^k / C^k$  and  $\hat{\mathbf{W}}_t = \sum_{k=1}^K W_t^k / C^k$ . Recursion of the queue dynamics (27) and letting  $T \rightarrow \infty$  where  $\mathbf{W}_T \in \mathcal{R}_s$ , gives the standard expression for the infinite-horizon sum workload [26]:

$$\hat{\mathbf{W}}_0 := \hat{G}(\mathbf{a}) = \sup_{t \in \mathbb{N}} \hat{\mathbf{a}}(0, t] - (t - 1), \quad (28)$$

where  $\hat{G}$  represents the infinite-horizon normalized sum workload mapping; in other words, for all  $s$ ,  $\hat{\mathbf{W}}_s = \hat{G}(\mathbf{a}|_{(s, \infty]})$  represent the normalized sum workload at time  $s$ , under arrival sequence  $\mathbf{a}$ .

To prove the lemma we use the fact that the rate region  $\mathcal{R}_s$  is simplex, hence  $\hat{\mathbf{W}}_s \leq 1 \Leftrightarrow \mathbf{W}_s \in \mathcal{R}_s$ . Thus, it suffices to show that there is a finite  $n'_0$  and a finite  $s$  such that  $\hat{G}(\mathbf{a}|_{(s, \infty]}) \leq 1$ , and for  $n \geq n'_0$ ,  $\hat{G}(\mathbf{a}^n|_{(s, \infty]}) \leq 1$ .

Since  $\mathbf{a} \in \mathcal{D}_{\boldsymbol{\mu}}^K$ , there is a  $t_0 < \infty$  such that for all  $\epsilon > 0$ ,  $t > t_0$  and  $k \in K$ ,  $\frac{a^k(0, t]}{t} \leq \mu^k + \epsilon C^k$ . Since  $\boldsymbol{\mu} \in \text{int}(\mathcal{R}_s)$ , we choose  $\epsilon = (1 - \hat{\boldsymbol{\mu}})/4K$ . We now have that for all  $t \geq t_0$ ,  $\frac{\hat{\mathbf{a}}(0, t]}{t} \leq \hat{\boldsymbol{\mu}} + \epsilon K = (1 + 3\hat{\boldsymbol{\mu}})/4 < 1$ . In other words, the workload at time zero is a function of only the arrivals within time  $(0, t_0]$  and hence,

$$\hat{\mathbf{W}}_0(\mathbf{a}) = \sup_{1 \leq t \leq t_0} \hat{\mathbf{a}}(0, t] - (t - 1). \quad (29)$$

Let  $s^* \leq t_0 < \infty$  be the minimum values of the optimizing  $t$ 's in the above equation. It can be shown [26, Lemma 5.4] that

$$\hat{G}(\mathbf{a}|_{(s^*, \infty]}) \leq 1.$$

It is known that  $\hat{G}$  is continuous on  $\mathcal{D}_{\hat{\boldsymbol{\mu}}}$  [24, Lemma 13] when  $\hat{\boldsymbol{\mu}} < 1$ . However, this together with continuity of shift mapping implies that for all  $\{\mathbf{a}^n\}$  such that  $\mathbf{a}^n$  converges to  $\mathbf{a}$  in scaled uniform topology,

$$\hat{G}(\mathbf{a}^n|_{(u, \infty]}) \rightarrow \hat{G}(\mathbf{a}|_{(u, \infty]}) \text{ for all } u \in [0, s^*]. \quad (30)$$

In particular, (30) implies that there exists  $n_0$  such that for all  $n \geq n_0$ , the normalized sum workload under arrival sequence  $\mathbf{a}^n$  at time  $u = s^*$  is no more than 1 packets, i.e.,

$$\hat{G}(\mathbf{a}^n|_{(s^*, \infty]}) < 1 \text{ for all } n \geq n_0.$$

However, since the rate region is a simplex, the workload vectors at time  $s^*$ , under  $\mathbf{a}$  and  $\mathbf{a}^n$  lie in the rate region  $\mathcal{R}_s$ . Hence, we have the assertion of the lemma.  $\blacksquare$

Next, we prove Lemma 6:

*Lemma 6:* Let  $\boldsymbol{\mu} \in \text{int}(\mathcal{R}_s)$ ,  $\mathbf{a}$  be an arrival sequence with rate  $\boldsymbol{\mu}$  with  $I^\sharp(\mathbf{a}) < +\infty$ , and  $\mathcal{W} = G^{\text{wc}}(\mathbf{a})$  be its corresponding steady state workload. For any  $\mathcal{W} \in {}^a G^{\text{wc}}$  there exists a sequence of arrivals  $\{\mathbf{a}^n \in \mathcal{D}_\mu^K\}$  such that  $\mathbf{a}^n$  converges to  $\mathbf{a}$  in the scaled uniform norm topology,  $G^{\text{wc}}(\mathbf{a}^n) \rightarrow G^{\text{wc}}(\mathbf{a})$ ,  $G^{\text{wc}}$  is continuous at  $\mathbf{a}^n$ , and  $I^\sharp(\mathbf{a}^n) \rightarrow I^\sharp(\mathbf{a})$ .

*Proof:* Lemma 5 implies that for any given arrival sequence  $\mathbf{a} \in \mathcal{D}_\mu^K$ , there exists a  $s^*$  such that

$$G^{\text{wc}}(\mathbf{a}) = G_{s^*}^{\text{wc}}(\mathbf{a}|_{(0,s^*]}).$$

However,  $G_{s^*}^{\text{wc}}$  is quasi-continuous on  $\mathbb{R}_+^{K \times s^*}$ . This implies that there exists a sequence of finite arrivals  $\{\hat{\mathbf{a}}^n|_{(0,s^*]}\}$  such that

- 1)  $\hat{\mathbf{a}}^n|_{(0,s^*]} \rightarrow \mathbf{a}|_{(0,s^*]}$ ;
- 2)  $G_{s^*}^{\text{wc}}$  is continuous at  $\hat{\mathbf{a}}^n|_{(0,s^*]}$ ; and
- 3)  $G_{s^*}^{\text{wc}}(\hat{\mathbf{a}}^n|_{(0,s^*]}) \rightarrow G_{s^*}^{\text{wc}}(\mathbf{a}|_{(0,s^*]})$ .

Now construct the sequence of arrivals  $\{\mathbf{a}^n\}$  via concatenation of  $\mathbf{a}|_{(s^*,\infty]}$  and  $\hat{\mathbf{a}}^n|_{(0,s^*]}$ . It is immediate that  $\mathbf{a}^n \rightarrow \mathbf{a}$ .

Appealing to Lemma 5, for  $n$  large enough (greater than  $n_0$ ) we have

- 1) convergence of  $G^{\text{wc}}(\mathbf{a}^n)$  to  $G^{\text{wc}}(\mathbf{a})$  since

$$G^{\text{wc}}(\mathbf{a}^n) = G_{s^*}^{\text{wc}}(\mathbf{a}^n|_{(0,s^*]}) \rightarrow G_{s^*}^{\text{wc}}(\mathbf{a}|_{(0,s^*]}) = G^{\text{wc}}(\mathbf{a});$$

- 2) continuity of  $G^{\text{wc}}(\cdot)$  at  $\mathbf{a}^n$ . For any sequence converging to  $\mathbf{a}^n$  in  $\mathcal{D}_\mu^K$  by appealing to Lemma 5 we know that far enough along every sequence only the arrivals in  $(0, s^*]$  matter. Now using the fact that projection is continuous on  $\mathcal{D}_\mu^K$ , we get the result from the continuity of  $G_{s^*}^{\text{wc}}$  at  $\mathbf{a}^n|_{(0,s^*]}$ .

This establishes the quasi continuity of function  $G^{\text{wc}}$ . Lastly, Assumptions 3 and 2 ensure that  $I^\sharp(\mathbf{a}^n) \rightarrow I^\sharp(\mathbf{a})$ . ■

Finally, we prove Lemma 7.

*Lemma 7:* If  $\boldsymbol{\mu} \in \text{int}(\mathcal{R}_s)$ , the mapping  $G^{\text{wc}}(\cdot)$  is almost compact on  $\mathcal{D}_\mu^K$  with respect to the scaled uniform norm topology.

*Proof:* This follows almost exactly along the same lines as the proof of Lemma 6. For any  $\mathbf{a}^n \rightarrow \mathbf{a}$  we proved the existence of a  $n_0$  such that for  $n \geq n_0$  such that the workload vectors

only depended on the arrivals within time  $(0, t_0]$ . Thus, the proof of almost compactness simply follows from Lemma 4. Note that we have used the fact that the projection operator is continuous on  $\mathcal{D}_\mu^K$ .  $\blacksquare$

## APPENDIX D

### PROOF OF LEMMA 10

Next we prove Lemma 10 which gives the bounds on  $I_t$ .

*Lemma 10:* For  $K = 2$ ,  $\mathbf{b} \in \mathbb{R}_+^2$ ,  $I_t(\mathbf{b})$  can be bounded as

$$I_t(\mathbf{b}) \geq \min_{u \in (0, t]} u \sum_{k=1}^K \Lambda_1^* \left( \frac{1}{u} \left( \text{Proj}_{\mathbb{X}(u, \mathbf{b})}(\mathbf{0}) \right)^k \right)$$

and when  $\mathbf{b} \notin [0, 1]^2$ ,

$$I_t(\mathbf{b}) \leq \min_{u \in (0, t]} u \sum_{k=1}^K \Lambda_1^* \left( \frac{1}{u} (b^k + (u-1)H(\mathbf{b})^k) \right),$$

where we recall that the convex set  $\mathbb{X}(u, \mathbf{b}) \subseteq \mathbb{R}_+^2$  is defined as

$$\mathbb{X}(u, \mathbf{b}) := \{\mathbf{b} + \mathbf{v} : \mathbf{v} \in \mathbb{R}_+^2 \text{ and } v^1 + v^2 = (u-1)\}.$$

*Proof:* Let  $\mathbf{b} \in \mathbb{R}_+^2$ , time  $u \in (0, t]$ , arrival path  $\mathbf{a}|_{(0, t]} \in \mathbb{A}(u, \mathbf{b})$ , and  $\mathbf{W}_i \in \mathbb{R}_+^2$  be the workload vector at time  $-i$  for  $i \in (0, u]$ . Assume without loss of generality that  $t > 1$  as it is easy to see that both bounds turn to be  $I_1(\mathbf{b})$ . Owing to this we can also assume that  $u > 1$  since the terms corresponding to  $u = 1$  in both bounds evaluate to  $I_1(\mathbf{b})$ .

We first show the lowerbound (16). As we have noted earlier that for  $\mathbf{a}|_{(0, t]} \in \mathbb{A}(u, \mathbf{b})$ , the  $[\cdot]^+$  function can be removed from the queue dynamics. Hence, we have  $\mathbf{a}(0, u] = \mathbf{b} + \sum_{i=1}^{u-1} H(\mathbf{W}_i)$ , where  $\mathbf{W}_u \in \mathcal{R}_s$  and  $\mathbf{W}_i \notin \mathcal{R}_s$  for all  $i \in (0, u-1]$ . Using this and the fact that  $H(\mathbf{W}_i) \in \{\mathbf{v} \in \mathbb{R}_+^2 : v^1 + v^2 = 1\}$ , for all  $i \in (0, u-1]$ , we have  $\mathbf{a}(0, u] \in \mathbb{X}(u, \mathbf{b})$  where  $\mathbb{X}(u, \mathbf{b})$  is defined above. Now, given any point  $\mathbf{d} \in \mathbb{X}(u, \mathbf{b})$ , the constant-speed linear path with increments of  $\mathbf{d}/u$  is the minimum-cost path among all the paths with the same destination (using Property 1). In addition, among all the paths to destinations in  $\mathbb{X}(u, \mathbf{b})$ , the closest constant-speed linear paths  $\mathbf{a}^*|_{(0, u]}$  to the equal line is the minimum-cost path (using Property 2). Since the closest point in  $\mathbb{X}(u, \mathbf{b})$  to the equal line is  $\text{Proj}_{\mathbb{X}(u, \mathbf{b})}(\mathbf{0})$ , we have  $\mathbf{a}^*|_{(0, u]} = (\mathbf{a}_i^* = \frac{1}{u} \text{Proj}_{\mathbb{X}(u, \mathbf{b})}(\mathbf{0}), i \in (0, u])$ .

Since the set of paths with destination in  $\mathbb{X}(u, \mathbf{b})$  includes all paths in  $\mathbb{A}(u, \mathbf{b})$ , from (7) we have the lowerbound (16):

$$\begin{aligned} I_t(\mathbf{b}) &= \min_{u \in (0, t]} \inf_{\mathbf{x} \in \mathbb{A}(u, \mathbf{b})} \sum_{k=1}^K \sum_{i=1}^u \Lambda_1^*(x_i^k) \geq \min_{u \in (0, t]} \inf_{\mathbf{x} \in \mathbb{R}_+^{Kt}: \mathbf{x}(0, u) \in \mathbb{X}(u, \mathbf{b})} I_t(\mathbf{x}) \\ &= \min_{u \in (0, t]} \inf_{\mathbf{x} \in \mathbb{R}_+^{Ku}: \mathbf{x} \in \mathbb{X}(u, \mathbf{b})} I_u(\mathbf{x}) = \min_{u \in (0, t]} u \sum_{k=1}^K \Lambda_1^* \left( \frac{1}{u} \left( \text{Proj}_{\mathbb{X}(u, \mathbf{b})}(\mathbf{0}) \right)^k \right). \end{aligned}$$

To show the upperbound (17), we only need to show that the constant-speed linear path  $\mathbf{a}|_{(0, u]} = (\mathbf{a}_i = \frac{1}{u}(\mathbf{b} + (u-1)H(\mathbf{b})), i \in (0, u])$ , is in  $\mathbb{A}(u, \mathbf{b})$ , when  $\mathbf{b} \notin [0, 1]^2$ . Without loss of generality, we consider only when  $b^1 \geq b^2$  and  $b^1 \geq 1$ . In this case, we set  $H(\mathbf{b}) = (1, 0)$  and the queue dynamics gives

$$\mathbf{W}_i = \frac{(u-i)}{u}(\mathbf{b} + (u-1)(1, 0)) - (u-1-i)(1, 0),$$

for all  $i \in (0, u-1]$ . Since  $b^1 \geq 1$ , we have  $W_i^1 \geq 1$  and  $W_i^1 \geq W_i^2$ , and hence we can once again set  $H(\mathbf{W}_i) = (1, 0)$  for all  $i \in (0, u-1]$ . Hence,  $\mathbf{a}|_{(0, u]} \in \mathbb{A}(u, \mathbf{b})$ . ■

## REFERENCES

- [1] L. Tassiulas and A. Ephremides, "Stability properties of constrained queueing systems and scheduling policies for maximum throughput in multi-hop radio networks," *IEEE Trans. Autom. Control*, vol. 37, no. 12, pp. 1936–1948, 1992.
- [2] N. McKeown, A. Mekkittikul, V. Anantharam, and J. Walrand, "Achieving 100% throughput in an input-queued switch," *IEEE Trans. Comm.*, vol. 47, no. 8, pp. 1260–1267, 1999.
- [3] J. G. Dai and B. Prabhakar, "The throughput of data switches with and without speed-up," in *Proc. of IEEE Infocom.*, pp. 556–564, 2000.
- [4] M. Armony and N. Bambos, "Queueing dynamics and maximal throughput scheduling in switched processing systems," *Queueing Systems*, vol. 44, no. 3, pp. 209–252, 2003.
- [5] M. Andrews, A. Stolyar, K. Kumaran, R. Vijayakumar, K. Ramanan, and P. Whiting, "Scheduling in a queueing system with asynchronously varying service rates," *Probab. Eng. Inf. Sci.*, vol. 18, pp. 191–217, 2004.
- [6] J. .A. van Mieghem, "Dynamic scheduling with convex delay costs: The generalized  $c\mu$  rule," *Ann. Appl. Prob.*, vol. 5, no. 3, pp. 809–833, 1995.
- [7] A. L. Stolyar, "MaxWeight scheduling in a generalized switch: State space collapse and equivalent workload minimization in heavy traffic," *Ann. Appl. Prob.*, vol. 14, no. 1, pp. 1–53, 2004.
- [8] A. Mandelbaum and A. L. Stolyar, "Scheduling flexible servers with convex delay costs: Heavy-traffic optimality of the generalized  $c\mu$ -rule," *Oper. Res.*, vol. 52, no. 6, pp. 836–855, 2004.
- [9] D. Bertsimas, I. Paschalidis, and J. Tsitsiklis, "Asymptotic buffer overflow probabilities in multiclass multiplexers: an optimal control approach," *IEEE Trans. Autom. Control*, vol. 43, no. 3, pp. 315–335, 1998.
- [10] A. L. Stolyar and K. Ramanan, "Largest weighted delay first scheduling: Large deviations and optimality," *Annals of Applied Probabilities*, vol. 11, no. 1, pp. 1–48, Feb. 2001.

- [11] A. Ephremides and B. Hajek, "Information theory and communication networks: An unconsummated union," *IEEE Trans. Info. Th.*, vol. 44, no. 6, pp. 2416–2434, 1998.
- [12] S. Shakkottai, "Effective Capacity and QoS for wireless scheduling," *IEEE Trans. Autom. Control*, vol. 53, no. 3, pp. 749–761, Apr. 2008.
- [13] L. Ying, R. Srikant, A. Eryilmaz, and G. Dullerud, "A Large Deviations analysis of scheduling in wireless networks," *IEEE Trans. Inf. Th.*, vol. 52, no. 11, pp. 5088–5098, Nov. 2006.
- [14] V. G. Subramanian, "Large deviations of max-weight scheduling policies of convex rate regions," in *Proc. ITA*, 2008.
- [15] V. G. Subramanian, "Large deviations of max-weight scheduling policies of convex rate regions," *preprint*.
- [16] C.-W. Yang, A. Wierman, S. Shakkottai, and M. Harchol-Balter, "Tail asymptotics for policies favoring short jobs in a many-flows regime," *SIGMETRICS Perform. Eval. Rev.*, vol. 34, no. 1, pp. 97–108, 2006.
- [17] C.-W. Yang, A. Wierman, S. Shakkottai, and M. Harchol-Balter, "Many flows asymptotics for SMART scheduling policies," *preprint*.
- [18] S. Shakkottai and R. Srikant, "Many-sources delay asymptotics with applications to priority queues," *Queueing Systems Theory and Applications (QUESTA)*, vol. 39, pp. 183–200, Oct. 2001.
- [19] A. Weiss, "A new technique for analyzing large traffic systems," *Advances in Applied Probability*, vol. 18, pp. 506–532, 1986.
- [20] D. D. Botvich and N. G. Duffield, "Large deviations, the shape of the loss curve, and economies of scale in large multiplexers," *Queueing System*, vol. 20, pp. 293–320, 1995.
- [21] C. Courcoubetis and R. Weber, "Buffer overflow asymptotics for a buffer handling many traffic sources," *Journal of Applied Probability*, vol. 33, pp. 886–903, 1996.
- [22] A. Simonian and J. Guibert, "Large deviations approximation for fluid queues fed by a large number of on/off sources," *IEEE JSAC*, vol. 13, no. 6, pp. 1017–1027, Aug. 1995.
- [23] D. J. Wischik, "The output of a switch, or, effective bandwidths for networks," *Queueing Systems Theory Appl.*, Vol. 32, No. 4, pp. 383–396, 1999.
- [24] D. J. Wischik, "Sample path large deviations for queues with many inputs," *Ann. Appl. Probab.*, 2001.
- [25] D. J. Wischik, "Moderate deviations in queueing theory," *preprint*.
- [26] A. Ganesh, N. O'Connell, and D. Wischik, *Big Queues*. Springer-Verlag, 2004.
- [27] E. Buffet and N. G. Duffield, "Exponential upper bounds via martingales for multiplexers with markovian arrivals," *J. Appl. Prob.*, vol. 31, pp. 1049–1060, 1994.
- [28] N. G. Duffield, "Exponential bounds for queues with markovian arrivals," *Queueing Systems*, vol. 17, pp. 413–430, 1994.
- [29] G. L. Choudhury, D. M. Lucantoni, and W. Whitt, "Squeezing the most out of atm," *IEEE Trans. Comm.*, vol. 44, pp. 203–217, Feb 1996.
- [30] A. Shwartz and A. Weiss, "Large deviations for performance analysis: Queues, communications, and computing," *Stochastic Modeling Series*, Chapman & Hall, London, 1995.
- [31] N. Likhanov and R. R. Mazumdar, "Cell loss asymptotics for buffers fed with a large number of independent stationary sources," *J. Appl. Probab.*, Vol. 36, No. 1, pp. 86–96, 1999.
- [32] M. Mandjes and J. H. Kim, "Large deviations for small buffers: An insensitivity result," *Queueing Syst.*, Vol. 37, No. 4, pp. 349–362, 2001.
- [33] M. Mandjes and S. Borst, "Overflow behavior in queues with many long-tailed inputs," *Adv. in Appl. Probab.*, Vol. 32, No. 4, pp.1150–1167, 2000.

- [34] C-W. Yang and S. Shakkottai, "Asymptotic evaluation of delay in the SRPT scheduler," *IEEE Trans. Automat. Control*, Vol. 51, No. 11, pp. 1848–1854, 2006.
- [35] O. Ozturk, R. R. Mazumdar and N. Likhanov, "Many sources asymptotics for networks with small buffers," *Queueing Syst.*, Vol. 46, No. 1–2, pp. 129–147, 2004.
- [36] S. Delas, R. R. Mazumdar and C. P. Rosenberg, "Tail asymptotics for HOL priority queues handling a large number of independent stationary sources," *Queueing Syst.*, Vol. 40, No. 2, pp. 183–204, 2002.
- [37] C. Kotopoulos and R. R. Mazumdar, "Buffer Occupancy and Delay Asymptotics in Multi-buffered Systems with Generalized Processor Sharing Handling a Large Number of Independent Traffic Streams," *preprint*.
- [38] M. Mandjes and M. van Uitert, "Sample-path large deviations for tandem and priority queues with Gaussian inputs," *Ann. Appl. Probab.*, Vol. 15, No. 2, pp. 1193–1226, 2005.
- [39] K. Debicki and M. Mandjes, "Exact overflow asymptotics for queues with many Gaussian inputs," Vol. 40, No. 3, pp. 704–720, 2003.
- [40] S. Kittipiyakul, P. Elia, and T. Javidi, "High-SNR analysis of outage-limited communications with bursty and delay-limited information," to appear in *IEEE Trans. Inf. Th.*
- [41] J. Garcia, "An extension of the Contraction Principle," *Journal of Theoretical Prob.*, vol. 17, no. 2, pp. 403–434, Apr. 2004.
- [42] J. Munkres, *Topology*, 2nd ed. Prentice Hall, 2000.
- [43] D. W. Muller, "Verteilungs-Invarianzprinzipien für das starke Gesetz der grossen Zahl," *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete*, Vol. 10, pp. 173–192, 1968.
- [44] W. Whitt, "Stochastic Abelian and Tauberian theorems," *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete*, Vol. 22, pp. 251–267, 1972.
- [45] A. A. Borovkov and A. I. Sahanenko, "Remarks on the convergence of random processes in nonseparable metric space and on the nonexistence of a Borel measure for processes in  $C(0, \infty)$ ," *Teor. Veroyatnost. i Primenen.*, Vol. 18, pp. 812–815, 1973.
- [46] A. A. Borovkov, "Convergence of distributions of functionals of random sequences and processes defined on the real line," *Proc. Steklov Inst. Math.*, Vol. 128, pp. 43–72, 1974.
- [47] A. I. Sahanenko, "The convergence of the distributions of functionals of processes that are defined on the whole axis," *Sibirsk. Mat. Ž.*, Vol. 15, pp. 102–119, 237, 1974.
- [48] H. Bauer, "On invariance principles with limit processes satisfying strong laws," *Z. Wahrsch. Verw. Gebiete*, Vol. 58, No. 2, pp. 257–265, 1981.
- [49] J. D. Deuschel and D. W. Stroock, *Large deviations*, Pure and Applied Mathematics Series, Vol. 137, Academic Press, 1989.
- [50] A. J. Ganesh and N. O'Connell, "A large deviations principle with queueing applications," *Stochastics and Stochastic Reports*, vol. 73, no. 1–2, pp. 25–35, 2002.
- [51] A. Dembo and O. Zeitouni, *Large Deviations techniques and applications*, 2nd ed. Springer, 1998.
- [52] S. Kittipiyakul, "Cross-layer optimization for transmission of delay-sensitive and bursty traffic in wireless systems," Ph.D. dissertation, University of California at San Diego, 2008.
- [53] R. J. R. Cruise, "A Scaling Framework for the Many Sources Asymptotic, through Large Deviations," Talk, *Young European Queueing Theorists Workshop II*, EURANDOM, Eindhoven, Dec 2008.
- [54] R. T. Rockafellar, *Convex Analysis*. Princeton Mathematical Series, No. 28, Princeton University, 1970.

- [55] A. W. Marshall and I. Olkin, *Inequalities: Theory of majorization and its applications*. Mathematics in Science and Engineering Series, Vol. 143, Academic Press Inc., 179.
- [56] R. T. Rockafellar and R. J.-B. Wets, *Variational Analysis*. Fundamental Principles of Mathematical Sciences Series, vol. 317, Springer-Verlag, 1998.
- [57] J.-P. Aubin and A. Cellina, *Differential Inclusions*. Fundamental Principles of Mathematical Sciences Series, vol. 264, Springer-Verlag, 1984.
- [58] M. Matejdes, “Sur les sélecteurs des multifonctions,” *Math. Slovaca*, vol. 37, no. 1, pp. 111–124, 1987.
- [59] M. Matejdes, “On the cliquish, quasicontinuous and measurable selections,” *Math. Bohem.*, vol. 116, no. 2, pp. 170–173, 1991.
- [60] J. Cao and W. B. Moors, “Quasicontinuous selections of upper continuous set-valued mapping,” *Real Anal. Exchange*, vol. 31, no. 1, pp. 63–71, 2005/2006.
- [61] R. Cazacu and J. D. Lawson, “Quasicontinuous functions, domains and extended calculus,” *Appl. Gen. Topol.*, vol. 8, no. 1, pp. 1–33, 2007.