

# On Characterizing the Delay-Performance of Wireless Scheduling Algorithms

Xiaojun Lin

Center for Wireless Systems and Applications  
School of Electrical and Computer Engineering, Purdue University  
West Lafayette, IN 47907, U.S.A.  
Email: linx@ecn.purdue.edu.

**Abstract**—In this paper we study the problem of characterizing the delay performance of wireless scheduling algorithms. In wireless networks operated under these wireless scheduling algorithms, there often exists a tight coupling between the service-rate process, the system backlog process, the arrival process and the channel variations. Although one can use sample-path large-deviation techniques to form an estimate of the delay-violation probability under a given offered load, the formulation leads to a multi-dimensional calculus of variations problem that is often very difficult to solve. In this paper, we present a new technique for addressing this complexity issue. Using ideas from the Lyapunov function approach in control theory, this technique maps the complex multi-dimensional calculus of variations problem to a one-dimensional calculus of variations problem, and the latter is often much easier to solve. We believe that this technique can potentially be used to study the delay-performance of a large class of wireless scheduling algorithms.

## I. INTRODUCTION

A wireless network may be modeled as a system of queues with time-varying service rates. The variability in service rates is due to a number of factors. First, channel fading and mobility can lead to variations in the link capacity even if the transmission power is fixed. Second, the transmission power can vary over time according to the power control policy. Third, due to radio interference, it is usually preferable to schedule only a subset of links to be active at each time, and to alternate the subset of activated links over time. All of these factors lead to a variable service rate at each link.

When one is to study the performance of any system that involves queues, the first question we can ask is whether the system is *stable* or not. Here, *stability* means that all queue length (or equivalently, the delay experienced by the packets) remains finite. Conversely, we can ask the question that, in order to maintain stability, what is the largest offered load that the system can carry. In other words, what is the capacity region of the system subject to stability. For wireless networks, these questions have led to results on *throughput-optimal* scheduling algorithms for scheduling wireless resources. (Here we use the term *scheduling* in the broader sense, i.e., it can include various control mechanisms at the MAC/PHY layer, e.g., link scheduling, power control, and adaptive coding/modulation.) A scheduling algorithm is *throughput-optimal* if this algorithm can sustain the largest offered load while keeping the system stable. In other

words, if this algorithm cannot stabilize the system, no other algorithms can. For example, one such throughput-optimal scheduling algorithm is the algorithm proposed in the seminal work by Tassiulas and Ephremides in [1]. This algorithm chooses at each time, among all possible schedules, the one that maximizes the sum of the queue-weighted-rate over all links. This algorithm has been shown to be throughput-optimal, and it has been the basis for many other throughput-optimal scheduling algorithms for both cellular and multihop wireless networks.

Once we know about stability, we are then tempted to ask further questions regarding the distribution of queue length (or delay). For example, at a given offered load, what is the probability that the delay experienced by a packet is greater than a given threshold? Or, conversely, what is the largest offered load that the system can support at a given delay constraint? (In other words, what is the *effective capacity region* of the system under delay constraints?) Clearly, these questions are important for applications that require more stringent delay guarantees than just stability.

These delay characterization problem for wireless networks can be difficult to solve. Here we draw a comparison to the delay characterization problem in wireline networks. In wireline networks, even through the exact delay distribution can be difficult to obtain, there have been a large body of work, especially those using large-deviation techniques, to obtain sharp estimates of the delay violation probability of a queue. These wireline network results usually assume that the service rate of the queue is fixed (i.e., time-invariant), and the packet arrival process is known. These results allow us to compute the *effective bandwidth* of the arrival process from its (known) statistics [2]–[7], which can then be used to determine the traffic carrying capability of the queue at a given delay constraint. In contrast, in wireless networks, the service rate is time-varying. If the service rate process is again known *a priori*, large-deviation techniques can be used to compute the *effective capacity* of the service rate process [8], [9], which is a notion similar to the *effective bandwidth* of the arrival process. This effective capacity can again be used to determine the traffic carrying capability of the queue at a given delay constraint. Unfortunately, under many wireless scheduling algorithms of interest, even the service rate process is unknown *a priori*. For example, for

a system operated under the throughput-optimal Tassioulas-Ephremides algorithm of [1], or any queue-length based scheduling algorithms, the service rates depend on the queue length, which in turn depend on the arrival process and the channel state, etc. Hence, the statistics of the service rate process is unknown before hand. In this case, the delay characterization problem is known to be very difficult. For these systems, although it is still possible to use sample-path large-deviation techniques to form an estimate of the delay-violation probability [10]–[12], such a formulation leads to a multi-dimensional calculus of variations problem. Due to the complex coupling between the service rate, the queue length, the arrival process, and the channel state, this multi-dimensional calculus of variations problem is very difficult to solve. Prior successes have been limited to *simple* systems: either the problem has some restrictive structure (e.g., symmetry among all links) [11], or the size of the system is very small (e.g., two links) [10], [12], [13].

In this paper, we present a new approach to address this complexity issue. Motivated by the Lyapunov function approach for proving stability of complex systems, we provide a technique that maps the complex multi-dimensional calculus of variations problem into a one-dimensional calculus of variations problem, and the latter is often very easy to solve. The solution to the one-dimensional calculus of variations problem will then provide us with an upper bound estimate of the delay violation probability, and consequently, a lower bound estimate of the effective capacity region of the system. For many practical applications, the resulting effective capacity region is useful because the delay constraint is known to be satisfied.

We believe that this marriage between sample-path large-deviations and Lyapunov functions can develop into a powerful theory to characterize the delay performance of wireless systems under sophisticated scheduling algorithms. We can potentially lower the difficulty level of the delay-characterization problem to that of a stability problem. In other words, for any scheduling algorithm that is provably stable, which usually means that there exists a Lyapunov function, we could then apply this theory to characterize the delay performance. We provide an example of how this approach can be used to solve a more difficult problem than those studied in the literature.

The rest of the paper is organized as follows. We present the network model in Section II. We review a formulation of the sample-path large-deviation principle, and identify the complexity of the associated calculus of variations problem in Section III. Then, in Section IV, we provide a Lyapunov function based approach to address the complexity issue. In Section V, we provide an example to show how such an approach can be used. Then we conclude.

## II. THE SYSTEM MODEL

We consider the following model for a wireless system with  $L$  links. In order to model channel fading, we assume that the system can be in one of  $S$  states. We assume a slotted system, and denote the state of the system at time  $t$  to be

$C(t)$ . Further, we assume that the states  $C(t), t = 1, 2, \dots$  are *i.i.d.*, and let  $p_j = \mathbf{P}[C(t) = j]$  denote the probability that the state of the system at time  $t$  is  $j$ . Let  $\vec{p} = [p_1, \dots, p_S]$ . For ease of exposition, in the rest of the paper we also define  $\Phi_j(t) = \mathbf{1}_{\{C(t)=j\}}$  to be the indicator function that the state of the system at time  $t$  is  $j$ . Let  $\vec{\Phi}(t) = [\Phi_1(t), \dots, \Phi_S(t)]$ . Clearly, there is a one-to-one mapping between  $C(t)$  and  $\vec{\Phi}(t)$ .

Each link corresponds to a queue with time-varying service rate. The arrivals at each link  $i$  are at a constant rate  $\lambda_i$ . The service offered to link  $i$  is determined by the scheduling algorithm, and in general correlates with the service at other links and depends on the system backlog. Let  $X_i(t)$  denote the backlog at link  $i$  at time  $t$ , and let  $\vec{X}(t) = [X_1(t), \dots, X_L(t)]$ . We assume that the service rate offered to link  $i$  is a function of the global backlog  $\vec{X}(t)$  and the system state  $C(t)$ . In particular, let  $D_{ij}(\vec{X})$  denote the service offered to link  $i$  when the state of the system is  $j$  and the global backlog is  $\vec{X}$ . The evolution of the backlog at link  $i$  is then given by

$$X_i(t+1) = [X_i(t) + \lambda_i - \sum_{j \in S} \Phi_j(t) D_{ij}(\vec{X}(t))]^+, \quad i = 1, \dots, L \quad (1)$$

where  $[\cdot]^+$  denotes the projection to  $[0, +\infty)$ .

Assume that the system is stationary and ergodic. In this paper, we will focus on studying the probability that the system backlog exceeds a certain threshold  $B$ . In particular, let  $\epsilon$  denote our target on the overflow probability, we would like to ensure that

$$\mathbf{P}[\|\vec{X}(0)\| \geq B] \leq \epsilon, \quad (2)$$

where  $\|\cdot\|$  is an appropriately chosen norm, and  $B$  is the overflow threshold. Note that the constraint in (2) is equivalent to a constraint on the delay-violation probability when the arrival rates  $\lambda_i$  are constant, because the two types of constraints are related by (see [9], [11])

$$\mathbf{P}[\text{Delay at link } i \geq d_i] = \mathbf{P}[X_i(0) \geq \lambda_i d_i].$$

Hence, in the rest of the paper we will often refer to (2) as a delay constraint.

Unfortunately, the problem of calculating the exact probability  $\mathbf{P}[\|\vec{X}(0)\| \geq B]$  is often mathematically intractable. In this paper, we are interested in using large-deviation techniques to compute estimates of this probability. We assume that the following large-deviation result for the backlog process  $\vec{X}(t)$  holds. That is, when  $B$  is large, the following limit exists

$$\lim_{B \rightarrow \infty} \frac{1}{B} \log \mathbf{P}[\|\vec{X}(0)\| \geq B] = -I_0(\vec{\lambda}), \quad (3)$$

where  $I_0(\vec{\lambda})$  can be determined from a sample-path large-deviation principle that we will describe in Section III. Equation (3) implies that, when  $B$  is large, the overflow probability can be approximated as

$$\mathbf{P}[\|\vec{X}(0)\| \geq B] \approx \exp(-BI_0(\vec{\lambda})).$$

Thus, the problem of estimating the overflow probability is reduced to that of computing the rate  $I_0(\vec{\lambda})$ . Alternatively, using the above approximation, in order to satisfy the constraint (2), we only need to ensure that

$$I_0(\vec{\lambda}) \geq \theta \triangleq -\frac{\log \epsilon}{B}. \quad (4)$$

We can then define the *effective capacity region* under the constraint (2) as the set of arrival rates  $\vec{\lambda}$  such that the above inequality holds.

### III. THE SAMPLE-PATH LARGE DEVIATION PRINCIPLE

In this paper, we will study the problem of computing the rate  $I_0(\vec{\lambda})$  and characterizing the effective capacity region under the constraint (2). We first describe how  $I_0(\vec{\lambda})$  can be determined from a sample-path large-deviation principle for the backlog process  $\vec{X}(t)$ . (Note that establishing such a large-deviation principle is not the main focus of the paper. We refer the readers to [10]–[12] for details on the technical assumptions under which such a large-deviation principle holds.)

#### A. Notations

We follow the convention in [10], [11]. For a large enough  $T$ , define the empirical measure process on the time interval  $[-T, 0]$  as

$$s_j^B(t) = \frac{1}{B} \sum_{l=0}^{B(T+t)} \mathbf{1}_{\{C(l)=j\}},$$

for  $t = \frac{k}{B} - T$ ,  $k = 0, \dots, BT$ , and by linear interpolation otherwise. Note that, in the above definition, we have scaled both the time and the magnitude. The quantity  $s_j^B(t)$  can be interpreted as the sum of the (scaled) time in  $[-T, t]$  that the system is at state  $j$ . Further, it is easy to check that  $\sum_{j \in \mathcal{S}} s_j^B(t) = t + T$  for all  $t \in [-T, 0]$ . Let  $\vec{s}^B(t) = [s_1^B(t), \dots, s_S^B(t)]$ . Further, let  $\phi_j^B(t) = \frac{d}{dt} s_j^B(t)$ . (Note that the derivative is well defined almost everywhere on  $[-T, 0]$  except when  $t = k/B - T$  for some integer  $k$ .) Let  $\vec{\phi}^B(t) = [\phi_1^B(t), \dots, \phi_S^B(t)]$ . Note that  $\sum_{j \in \mathcal{S}} \phi_j^B(t) = 1$  for almost all  $t$ .

Analogously, define the scaled backlog process as,

$$x_i^B(t) = \frac{1}{B} X_i(B(T+t)),$$

for  $t = \frac{k}{B} - T$ ,  $k = 0, \dots, BT$ , and by linear interpolation otherwise. Let  $\vec{x}^B(t) = [x_1^B(t), \dots, x_L^B(t)]$ . Note that according to (1), the backlog process  $\vec{x}^B(t)$  is related to the process  $\vec{\phi}^B(t)$  by

$$\begin{aligned} & \frac{x_i^B(t+1/B) - x_i^B(t)}{1/B} \\ &= \lambda_i - \sum_{j \in \mathcal{S}} D_{ij}(\vec{x}^B(t)) \int_t^{t+1/B} \phi_j^B(s) ds, \\ & \text{for } t = \frac{k}{B} - T, k = 0, \dots, BT. \end{aligned} \quad (5)$$

Thus, given a particular initial condition  $\vec{x}^B(-T)$ , Equation (5) defines a mapping  $\mathbf{f}^B$  from the empirical measure process

$\vec{s}^B(t)$  to the backlog process  $\vec{x}^B(t)$ . Further, although we have assumed  $\vec{s}^B(t)$  to be piecewise linear to begin with, the definition of the mapping  $\mathbf{f}^B$  can be naturally extended to all absolute continuous functions  $\vec{s}^B(t)$ .

#### B. The Large-Deviation Principle

Let  $B \rightarrow \infty$ . We now have a sequence of scaled random walks  $\vec{s}^B(t)$ , and they map to a sequence of scaled backlog processes  $\vec{x}^B(t)$  through the sequence of mappings  $\mathbf{f}^B$ . For any  $\vec{\phi} \geq 0$  and  $\sum_{j \in \mathcal{S}} \phi_j = 1$ , define  $H(\vec{\phi}|\vec{p}) = \sum_{j \in \mathcal{S}} \phi_j \log \frac{\phi_j}{p_j}$ . The sequence of empirical measure processes  $\vec{s}^B(t)$  are known to satisfy a sample-path large deviation principle [14, p176] with large-deviation rate-function  $I_s^T(\vec{s}(\cdot))$  given as follows:

$$I_s^T(\vec{s}(\cdot)) = \int_{-T}^0 H(\vec{\phi}(t)|\vec{p}) dt,$$

if  $\vec{s}(t)$  is absolute continuous and component-wise non-decreasing on  $[-T, 0]$ ,  $\vec{s}(-T) = 0$ , and  $\sum_{j \in \mathcal{S}} s_j(t) = t + T$  for all  $t$ ; where  $\vec{\phi}(t) = \frac{d}{dt} \vec{x}(t)$ . (Note that  $\vec{\phi}(t)$  is well defined almost everywhere on  $[-T, 0]$  since  $\vec{s}(t)$  is absolute continuous on  $[-T, 0]$ .) Otherwise,

$$I_s^T(\vec{s}(\cdot)) = +\infty.$$

Such a large-deviation principle means that, for any set  $\Gamma$  of trajectories on  $[-T, 0]$  that is a *continuity set* [14, p5] according to the *essential supremum norm* [14, p176, p352], the probability that the sequence of empirical measure processes  $\vec{s}^B(t)$  fall into  $\Gamma$  must satisfy

$$\lim_{B \rightarrow \infty} \frac{1}{B} \log \mathbf{P}[\vec{s}^B(\cdot) \in \Gamma] = - \inf_{\vec{s}(\cdot) \in \Gamma} I_s^T(\vec{s}(\cdot)). \quad (6)$$

The large-deviation rate-function  $I_s^T(\cdot)$  characterizes how rarely each trajectory is. Note that  $I_s^T(\vec{s}(\cdot)) \geq 0$  for all trajectory  $\vec{s}(\cdot)$ . The larger the value of  $I_s^T(\vec{s}(\cdot))$  is, the further the “empirical probability distribution”  $\vec{\phi}(t)$  deviates from the “prior probability distribution”  $\vec{p}$ . Hence, the less likely the trajectory  $\vec{s}(\cdot)$  will occur. Equation (6) reflects the well-known large-deviation philosophy that “rare events occur in the most-likely way.” Precisely, when  $B$  is large, the probability that the empirical measure process  $\vec{s}^B(t)$  falls into a set  $\Gamma$  is determined by the trajectory in  $\Gamma$  that is most likely to occur, i.e., with the smallest  $I_s^T(\vec{s}(\cdot))$ .

Next, assume that the sequence of mappings  $\mathbf{f}^B$  has a limiting mapping  $f$  that also maps any absolute continuous empirical measure process  $\vec{s}(t)$  to a backlog processes  $\vec{x}(t)$ . Assume that the limiting mapping  $f$  is of the form

$$\frac{d}{dt} x_i(t) = \lambda_i - \sum_{j \in \mathcal{S}} \phi_j(t) d_{ij}(\vec{x}(t)), \quad (7)$$

where  $\vec{\phi}(t) = \frac{d}{dt} \vec{s}(t)$ . (Note that this equation may be viewed as the limit of (5) when  $B \rightarrow \infty$ , although the function  $d_{ij}(\vec{x}(t))$  may not be exactly the same as  $D_{ij}(\vec{x}^B(t))$  as we will see in the example in Section V.) Further, assume that the sequence of mappings  $\mathbf{f}^B$  are *exponentially equivalent* to  $f$  [14, p130], and the mapping  $f$  is continuous (see [10] for

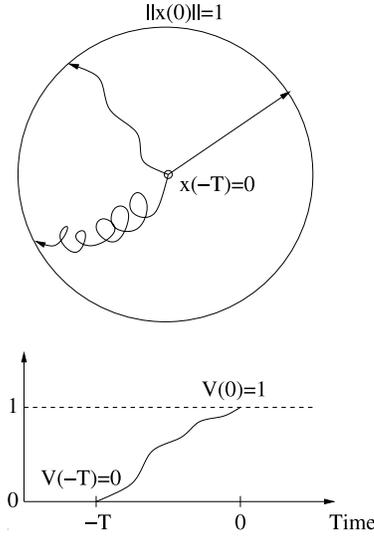


Fig. 1. Top: The overflow probability  $\mathbf{P}\{||\vec{X}(0)|| \geq B\}$  is related to the most likely path to overflow. Bottom: The technique that we present in Section IV maps any multi-dimensional path  $\vec{x}(t)$  to a one-dimensional path  $V(t)$ .

details of how the continuity of  $f$  may be verified). For any sequence of backlog processes that start from  $\vec{x}^B(-T) = 0$ , we can then invoke the contraction principle [14, p131] and obtain a sample-path large-deviation principle for the sequence of backlog processes  $\vec{x}^B(t)$  with large-deviation rate-function given by:

$$I_x^T(\vec{x}(\cdot)) = \inf_{\vec{s}(\cdot): \vec{x}(\cdot) = f(\vec{s}(\cdot))} \left\{ \int_{-T}^0 H(\vec{\phi}(t) | \vec{p}) dt \right\}$$

where  $\vec{\phi}(t) = \frac{d}{dt} \vec{s}(t)$ , and the infimum is taken over all empirical measure processes  $\vec{s}(\cdot)$  that map to the backlog process  $\vec{x}(\cdot)$  given that  $\vec{x}(-T) = 0$ , under the mapping  $f$ . Finally, the event of overflow corresponds to  $||\vec{x}^B(0)|| \geq 1$ . As  $B \rightarrow \infty$ , we have,

$$\begin{aligned} I_0(\vec{\lambda}) &\triangleq - \lim_{B \rightarrow \infty} \frac{1}{B} \log \mathbf{P}\{||\vec{x}^B(0)|| \geq 1\} \\ &= \inf \{ I_x^T(\vec{x}(\cdot)) \text{ over all trajectory } \vec{x}(\cdot) \text{ that} \\ &\quad \text{goes from } \vec{x}(-T) = 0 \text{ for some } T > 0 \\ &\quad \text{to } ||\vec{x}(0)|| = 1 \}. \end{aligned} \quad (8)$$

The trajectory that attains the infimum in (8) is often called the *most likely path to overflow* (see Figure 1). Clearly, in order to estimate the overflow probability  $\mathbf{P}\{||\vec{X}(0)|| \geq B\}$ , all we need is to find out which trajectory in (8) is the most-likely path to overflow.

### C. The “Path-Explosion” Challenge

Unfortunately, the infimum in (8) (a calculus of variations problem) is often very difficult to evaluate, because it is taken over an infinite number of multi-dimensional paths  $\vec{x}(t)$ . To see this, let us define the *local rate-function*

$$\begin{aligned} l(\vec{x}, \vec{y}) &= \inf \{ H(\vec{\phi} | \vec{p}) \text{ over all } \vec{\phi} \text{ such that } \sum_{j \in \mathcal{S}} \phi_j = 1 \\ &\quad \text{and } y_i = \lambda_i - \sum_{j \in \mathcal{S}} \phi_j d_{ij}(\vec{x}) \text{ for all } i \}. \end{aligned} \quad (9)$$

Then, we have

$$I_x^T(\vec{x}(\cdot)) = \int_{-T}^0 l(\vec{x}(t), \vec{x}'(t)) dt.$$

Note that the local rate-function  $l(\vec{x}, \vec{y})$  characterizes how *rarely* that, given  $\vec{x}(t) = \vec{x}$  at some time  $t$ ,  $\vec{x}(t)$  will follow the direction  $\frac{d}{dt} \vec{x}(t) = \vec{y}$  immediately after  $t$ . Suppose now we enforce a delay constraint in the form of  $\mathbf{P}\{||\vec{X}(0)|| \geq B\} \leq \epsilon$ . Using (8) to approximate this probability, we then need to ensure that

$$I_x^T(\vec{x}(\cdot)) = \int_{-T}^0 l(\vec{x}(t), \vec{x}'(t)) dt \geq \theta \triangleq -\log \epsilon / B \quad (10)$$

for *all* sample paths  $\vec{x}(t)$  that go from 0 at some past time  $-T$  to  $||\vec{x}(0)|| = 1$ . For advanced wireless scheduling algorithms like the Tassiulas-Ephremides algorithm [1], the complexity of enumerating all such paths soon becomes prohibitive. Prior successes have been limited to *simple* systems: either the problem has some restrictive structure (e.g., symmetry among all links) [11], or the size of the system is very small (e.g., two links) [10], [12], [13].

## IV. A NEW APPROACH COMBINING LARGE-DEVIATIONS WITH LYAPUNOV STABILITY

In this section, we will develop a general approach for solving problems (8) and (10). As we discussed earlier, finding the most likely path to overflow is often a very difficult problem. In this section, instead of solving the calculus of variations problem on the right-hand-side of (8), we construct a lower bound for it. That is, we will find a quantity  $\theta_0$  such that  $I_x^T(\vec{x}(\cdot)) \geq \theta_0$  for all trajectory  $\vec{x}(\cdot)$  that goes from  $\vec{x}(-T) = 0$  at some past time  $-T$  to  $||\vec{x}(0)|| = 1$ . Hence, we obtain an upper bound on the overflow probability. If  $\theta_0 \geq \theta$ , we then obtain a sufficient condition for meeting the constraint (4) on the overflow probability. Consequently, we also obtain a lower bound on the effective capacity region of the system under the constraint (2).

How to find such a lower bound  $\theta_0$ ? In this section, we present a technique that is motivated by the Lyapunov function approach for proving stability for complex systems [15]. Note that for a complex system like (1), it becomes difficult to even establish *stability*, i.e., to show that all queues will remain finite. To see this, take the limit  $B \rightarrow +\infty$  again for (5). The *fluid limit* of the system is governed by [16]

$$\frac{d}{dt} x_i(t) = \lambda_i - \sum_{j \in \mathcal{S}} p_j D_{ij}(\vec{x}(t)) \triangleq h_i(\vec{x}(t)), \text{ for all } i \quad (11)$$

or, in vector form

$$\frac{d}{dt} \vec{x}(t) = \vec{h}(\vec{x}(t)).$$

This fluid limit dynamics can be viewed as the *mean* behavior of the system. The original system would be stable if the solution of the ODE (ordinary differential equation) in (11) can be shown to converge to zero from any initial condition [16]. However, solving the above ordinary differential equation is often very difficult. Thus, it is usually impossible

to establish the stability of the system by directly solving the ODE. To circumvent this difficulty, we usually find a Lyapunov function  $V(\vec{x})$ , such that  $V(\vec{x}) \geq 0$ , and  $V(\vec{x}) = 0$  if and only if  $\vec{x} = 0$ . We then prove stability by showing a *negative drift* for  $V(\cdot)$ , i.e.

$$\frac{d}{dt}V(\vec{x}(t)) = \left(\frac{\partial V}{\partial \vec{x}}\right)^T \frac{d\vec{x}}{dt} \leq -\delta V(\vec{x}(t)), \quad (12)$$

where  $\delta$  is a small positive constant. Thus, if  $\vec{x}(t)$ , or equivalently,  $V(\vec{x}(t))$ , is away from zero, the negative drift will pull them back to zero. The negative drift then provides a sufficient condition for  $V(\vec{x}(t)) \rightarrow 0$ , which implies that  $\vec{x}(t) \rightarrow 0$ , as  $t \rightarrow \infty$ . Here, *the key to the Lyapunov function approach is to map the convergence of a multi-dimensional path  $\vec{x}(t)$ , to the convergence of a one-dimensional path  $V(\vec{x}(t))$ , which is then much easier to show.* Since  $V(\vec{x}(t)) \rightarrow 0$  provides a sufficient condition for all solutions  $\vec{x}(t)$  of the ODE to go to zero, the Lyapunov function approach allows us to identify a lower bound on the capacity region of the system (subject to stability).

Can we use a similar Lyapunov function approach to characterize the delay performance (and the overflow probability) of wireless scheduling algorithms? Indeed, Lyapunov functions have been used to solve other calculus of variations problems in the control literature. We next demonstrate how such an approach can be used for the delay-characterization problem. Without loss of generality, assume that the Lyapunov function  $V(\cdot)$  are chosen such that  $\|\vec{x}\| = 1$  implies  $V(\vec{x}) \geq 1$ . According to (12), for any  $w > -\delta V(\vec{x}(t))$ , the trajectory with  $\frac{d}{dt}V(\vec{x}(t)) = w$  becomes a ‘‘rare’’ event. Define

$$l_V(v, w) = \inf\{l(\vec{x}, \vec{y}) \mid \text{over all } (\vec{x}, \vec{y}) \text{ such that } V(\vec{x}) = v \text{ and } \left(\frac{\partial V}{\partial \vec{x}}\right)^T \vec{y} = w\}. \quad (13)$$

Let us abuse notation and let  $V(t) = V(\vec{x}(t))$ . Compared with (9),  $l_V(v, w)$  becomes the local rate-function for  $V(t)$ , i.e., it characterizes how *rarely* that, given  $V(t) = v$  at some time  $t$ ,  $V(t)$  will follow the direction  $\frac{d}{dt}V(t) = w$  immediately after  $t$ . It is easy to show that  $I_x^T(\vec{x}(\cdot))$  in (10) satisfies

$$I_x^T(\vec{x}(\cdot)) = \int_{-T}^0 l(\vec{x}(t), \vec{x}'(t))dt \geq \int_{-T}^0 l_V(V(t), V'(t))dt.$$

Let

$$\theta_0 = \inf\left\{\int_{-T}^0 l_V(V(t), V'(t))dt \mid \text{over all trajectory } V(\cdot) \text{ that goes from } V(-T) = 0 \text{ for some } T > 0 \text{ to } V(0) = 1\right\}. \quad (14)$$

We then obtain a lower bound  $\theta_0$  for the calculus of variations problem (8). It is also easy to see that a sufficient condition for all samples paths  $\vec{x}(t)$  to meet the constraint (10) is

$$\int_{-T}^0 l_V(V(t), V'(t))dt \geq \theta \quad (15)$$

for all *one-dimension* path  $V(t)$  that goes from  $V(-T) = 0$  to  $V(0) = 1$ . Again, we have successfully reduced the original multi-dimensional calculus of variations problem to a one-dimensional problem. The one-dimensional calculus of variations problem in (14) and (15) is usually much easier to solve (Fig. 1).

*Remark:* Lyapunov functions have been used in the control literature to solve other calculus of variations problems. Often, the key to success of such an approach is to find the right Lyapunov function. The unique feature of the scheduling problem studied in this paper is that the Lyapunov function for stability automatically becomes the suitable Lyapunov function for the calculus of variations problem. Note that for any scheduling algorithm that is provably stable, which usually means that there exists a Lyapunov function for stability, we can then apply the above techniques to characterize the delay performance. In other words, the difficulty level of the delay-characterization problem is reduced to that of a stability problem. Since (15) is a sufficient condition to (10), we can obtain an upper bound on the overflow probability, and correspondingly, if a constraint on the overflow probability is imposed, we obtain a lower bound on the effective capacity region. The hope of this approach is that, if the function  $V(\cdot)$  is appropriately chosen, we may recover a large fraction of, or even the entire effective capacity region.

## V. AN EXAMPLE

In this section, we apply the methodology of Section IV to the delay-characterization problem in [11]. The model of [11] is a base-station serving  $N$  users (Fig. 2). Packets for user  $i$  arrive at a constant rate  $\lambda_i$ . Only one user can be scheduled for transmission at any time. The fading channel between the base-station and each user is *i.i.d.*. At each time-slot, a user’s channel is ON with probability  $p$ , and OFF with probability  $1-p$ . Let  $F$  denote the bandwidth of the system. Hence, if a user’s channel is ON and it is scheduled for transmission, its service rate is  $F$ . The throughput-optimal Tassioulas-Ephremides algorithm [1] in this case is the QLB (Queue-Length Based) algorithm, i.e., the base-station should schedule the ON user with the longest queue [11]. The more challenging question is to determine the effective capacity region of the system, subject to the buffer overflow constraint  $\mathbf{P}[\max_{i=1, \dots, N} X_i \geq B] \leq \epsilon$ , where  $X_i$  is the random variable that denotes the backlog of user  $i$ . The authors of [11] assume that all users have the same offered load, i.e.,  $\lambda_i = \lambda$  for all  $i = 1, 2, \dots, N$ . Under this assumption, the most likely path to overflow in (8) can be explicitly solved. They then establish the following effective capacity region:

$$N\lambda \leq \min_{1 \leq M \leq N} \frac{N}{\theta} \log \left[ (1-p)^M + (1 - (1-p)^M) \exp\left(-\frac{F\theta}{M}\right) \right] \quad (16)$$

where  $\theta = -\log \epsilon / B$ . However, for non-identical offered loads, it appears very difficult to follow the solution approach of [11].

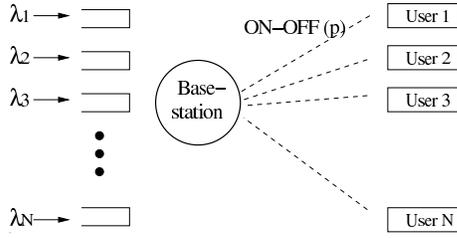


Fig. 2. The scheduling problem in cellular networks under fading channel.

We now use the methodology of Section IV to solve the delay characterization problem when the offered loads  $\lambda_i$  are non-identical. Following the notations in Sections II and III, the set of possible channel states are  $\mathcal{S} = \{(a_1, \dots, a_N) | a_1, \dots, a_N = \text{ON or OFF}\}$ . The probability that the channel state  $C(t)$  at time  $t$  is  $j$  is given by,

$$p_j = p^{n(j)}(1-p)^{N-n(j)},$$

where  $n(j)$  is the number of users with ON channel at state  $j$ . When the state is  $j$  and the system backlog is  $\vec{x}$ , let  $\mathcal{I}_1(\vec{x}, j)$  denote the set of those users whose channels are ON and who have the (identically) largest queue  $x_i$  among ON users. The evolution of the backlog is then given by (1), where the function can be chosen as  $D_{ij}(\vec{x}(t)) = F/|\mathcal{I}_1(\vec{x}(t), j)|$  if  $i \in \mathcal{I}_1(\vec{x}(t), j)$ , and  $D_{ij}(\vec{x}(t)) = 0$  otherwise. We can define  $\mathbf{f}^B$  according to (5). As  $B \rightarrow \infty$ , the limiting mapping  $f$  is given by (7). It is easy to show that the function  $d_{ij}(\cdot)$  in (7) must satisfy  $d_{ij}(\vec{x}(t)) = 0$  if  $i \notin \mathcal{I}_1(\vec{x}(t), j)$ . Further, if the set  $\mathcal{I}_1(\vec{x}(t), j)$  contains only one user  $i$ , i.e., there is a unique ON user  $i$  that has the largest queue at time  $t$ , then  $d_{ij}(\vec{x}(t)) = F$ . However, if  $\mathcal{I}_1(\vec{x}(t), j)$  contains multiple users, the definition of  $d_{ij}(\vec{x}(t))$  becomes somewhat involved [10]–[12]. Roughly speaking,  $d_{ij}(\vec{x}(t))$  should be defined so that the users in  $\mathcal{I}_1(\vec{x}(t), j)$  can maintain identical queues as much as possible. Regardless of the exact form of  $d_{ij}(\cdot)$ , the following relationship can be shown. For a given trajectory  $\vec{x}(\cdot)$ , let  $\mathcal{I}_1(\vec{x}(t)) = \{i | x_i(t) = \max_k x_k(t)\}$  be the set of users with the (identically) largest queue at time  $t$ . (Note that  $\mathcal{I}_1(\vec{x})$  is different from  $\mathcal{I}_1(\vec{x}, j)$  since in the definition of  $\mathcal{I}_1(\vec{x})$  we do not check whether a user is ON or OFF.) Further, let  $\mathcal{I}_2(\vec{x}(t), \vec{x}'(t)) = \{i \in \mathcal{I}_1(\vec{x}(t)) | \frac{d}{dt}x_i(t) = \max_{k \in \mathcal{I}_1(\vec{x}(t))} \frac{d}{dt}x_k(t)\}$ . That is,  $\mathcal{I}_2(\vec{x}(t), \vec{x}'(t))$  is the set of users that, among those users with the largest queue at time  $t$ , also have the largest queue growth rate. In other words, these set of users will have the largest queue *immediately after time*  $t$ . Then, immediately after time  $t$ , as long as one user in  $\mathcal{I}_2(\vec{x}(t), \vec{x}'(t))$  is ON, this group of users collectively must receive the full service rate  $F$ . Therefore, using (7), we must have

$$\begin{aligned} & \sum_{i \in \mathcal{I}_2(\vec{x}(t), \vec{x}'(t))} \frac{d}{dt}x_i(t) \\ &= \sum_{i \in \mathcal{I}_2(\vec{x}(t), \vec{x}'(t))} \lambda_i - F \sum_{j \in \mathcal{S}(\mathcal{I}_2(\vec{x}(t), \vec{x}'(t)))} \phi_j, \end{aligned} \quad (17)$$

where for any subset  $\mathcal{A} \subset \{1, 2, \dots, N\}$ ,  $\mathcal{S}(\mathcal{A})$  denotes the set of states  $j$  such that some user  $i \in \mathcal{A}$  is ON.

We can then use (9) and write down the formulation for  $l(\vec{x}, \vec{y})$  as

$$\begin{aligned} l(\vec{x}, \vec{y}) &= \inf_{\sum_{j \in \mathcal{S}} \phi_j = 1} H(\vec{\phi} | \vec{p}) \\ &\text{subject to} \quad y_i = \lambda_i - \sum_{j \in \mathcal{S}} \phi_j d_{ij}(\vec{x}) \text{ for all } i. \end{aligned}$$

Let the Lyapunov function be  $V(\vec{x}) = \max_i x_i$ . Although  $V(\cdot)$  is not differentiable at every point, it is sufficient to deal with its one-sided derivative. In the rest of this section, when we use the notation  $\frac{dg(t)}{dt}$ , we will mean  $\lim_{s \downarrow 0} \frac{g(t+s) - g(t)}{s}$ . We thus have

$$\frac{dV(\vec{x}(t))}{dt} = \max_{i \in \mathcal{I}_1(\vec{x}(t))} \frac{dx_i(t)}{dt},$$

(Recall that  $\mathcal{I}_1(\vec{x}) = \{i | x_i = \max_k x_k\}$  is the set of users with the identically largest queue when the system backlog is  $\vec{x}$ .) Thus, using (13), we have,

$$\begin{aligned} l_V(v, w) &= \inf_{\vec{x}, \vec{y}} l(\vec{x}, \vec{y}) \\ &\text{subject to} \quad \max_i x_i = v \\ &\quad \max_{i \in \mathcal{I}_1(\vec{x})} y_i = w. \end{aligned}$$

Combining the above two optimization problems, we thus have

$$\begin{aligned} l_V(v, w) &= \inf_{\sum_{j \in \mathcal{S}} \phi_j = 1} H(\vec{\phi} | \vec{p}) \\ &\text{subject to} \quad \max_i x_i = v \\ &\quad \max_{i \in \mathcal{I}_1(\vec{x})} y_i = w \\ &\quad y_i = \lambda_i - \sum_{j \in \mathcal{S}} \phi_j d_{ij}(\vec{x}), \text{ for all } i. \end{aligned} \quad (18)$$

This optimization problem is still quite difficult to solve. We will be contented to obtain a lower bound for the optimal value. First, recall that

$$\mathcal{I}_2(\vec{x}, \vec{y}) = \{i \in \mathcal{I}_1(\vec{x}) | y_i = \max_{k \in \mathcal{I}_1(\vec{x})} y_k\}.$$

Let us first compute the infimum of (18) for all  $\vec{x}, \vec{y}$  such that  $\mathcal{I}_2(\vec{x}, \vec{y}) = \mathcal{M}$ , where  $\mathcal{M}$  is a given subset of  $\{1, \dots, N\}$ . This sub-optimization can be written as

$$\begin{aligned} l_{V, \mathcal{M}}(v, w) &= \inf_{\sum_{j \in \mathcal{S}} \phi_j = 1} H(\vec{\phi} | \vec{p}) \\ &\text{subject to} \quad x_i = v \text{ for } i \in \mathcal{M} \\ &\quad x_i \leq v \text{ for } i \notin \mathcal{M} \\ &\quad y_i = w \text{ for } i \in \mathcal{M} \\ &\quad y_i < w \text{ for } i \in \mathcal{I}_1(\vec{x}) / \mathcal{M} \\ &\quad y_i = \lambda_i - \sum_{j \in \mathcal{S}} \phi_j d_{ij}(\vec{x}) \text{ for all } i. \end{aligned} \quad (19)$$

Note that for all  $i \in \mathcal{M} = \mathcal{I}_2(\vec{x}, \vec{y})$ , we have

$$w = \lambda_i - \sum_{j \in \mathcal{S}} \phi_j d_{ij}(\vec{x}).$$

Summing over all  $i \in \mathcal{I}_2(\vec{x}, \vec{y})$ , and using (17), we have,

$$|\mathcal{M}|w = \sum_{i \in \mathcal{M}} \lambda_i - F \sum_{j \in \mathcal{S}(\mathcal{M})} \phi_j,$$

where  $|\mathcal{M}|$  denotes the cardinality of  $\mathcal{M}$ . (Recall also that  $\mathcal{S}(\mathcal{M})$  is the set of states  $j$  such that some user in  $\mathcal{M}$  is ON.) We can then relax the constraints of (19) as:

$$\begin{aligned} \tilde{l}_{V, \mathcal{M}}(v, w) &= \inf_{\substack{\phi_j=1 \\ j \in \mathcal{S}}} H(\vec{\phi} | \vec{p}) & (20) \\ \text{subject to } |\mathcal{M}|w &= \sum_{i \in \mathcal{M}} \lambda_i - F \sum_{j \in \mathcal{S}(\mathcal{M})} \phi_j. \end{aligned}$$

This subproblem is solved in [11], and the solution is given by

$$\tilde{l}_{V, \mathcal{M}}(v, w) = u \log \frac{u}{1 - (1-p)^{|\mathcal{M}|}} + (1-u) \log \frac{1-u}{(1-p)^{|\mathcal{M}|}}.$$

where  $u = \frac{\sum_{i \in \mathcal{M}} \lambda_i - |\mathcal{M}|w}{F}$ . Let

$$\tilde{l}_V(v, w) = \inf_{\mathcal{M} \subset \{1, 2, \dots, N\}} \tilde{l}_{V, \mathcal{M}}(v, w). \quad (21)$$

Since (20) is a relaxation of (19), we have,

$$l_V(v, w) = \inf_{\mathcal{M}} l_{V, \mathcal{M}}(v, w) \geq \inf_{\mathcal{M}} \tilde{l}_{V, \mathcal{M}}(v, w) = \tilde{l}_V(v, w).$$

Therefore, in order to ensure that

$$\int_{-T}^0 l(\vec{x}(t), \vec{x}'(t)) dt \geq \theta$$

for all trajectory  $\vec{x}(\cdot)$  that goes from  $\vec{x}(-T) = 0$  at some past time  $-T$  to  $\|\vec{x}(0)\| = 1$ , it is sufficient to ensure that

$$\begin{aligned} \theta &\leq \inf \left\{ \int_{-T}^0 \tilde{l}_V(V(t), V'(t)) dt \mid \text{over all trajectory} \right. \\ &\quad \left. V(\cdot) \text{ that goes from } V(-T) = 0 \text{ for some } T > 0 \right. \\ &\quad \left. \text{to } V(0) = 1 \right\}. \end{aligned} \quad (22)$$

Note that  $\tilde{l}_V(v, w)$  in (21) does not depend on  $v$ . Therefore, the trajectory  $V(\cdot)$  that attains the infimum in (22) is in fact very easy to solve [17, p520], and the infimum is equal to  $\inf_{w \geq 0} \tilde{l}_V(v, w)/w$ . Therefore, using the definition of  $\tilde{l}_V(\cdot, \cdot)$  and  $\tilde{l}_{V, \mathcal{M}}(\cdot, \cdot)$ , it is then sufficient to ensure that

$$\begin{aligned} \theta &\leq \inf_{w \geq 0} \frac{1}{w} \tilde{l}_V(v, w) \\ &= \inf_{\mathcal{M}} \inf_{w \geq 0} \frac{1}{w} \tilde{l}_{V, \mathcal{M}}(v, w) \\ &= \inf_{\mathcal{M}} \inf_{0 \leq u \leq \frac{\sum_{i \in \mathcal{M}} \lambda_i}{F}} \frac{|\mathcal{M}|}{\sum_{i \in \mathcal{M}} \lambda_i - uF} D_{|\mathcal{M}|}(u | p) \end{aligned}$$

where

$$\begin{aligned} D_{|\mathcal{M}|}(u | p) &= \\ &= u \log \frac{u}{1 - (1-p)^{|\mathcal{M}|}} + (1-u) \log \frac{1-u}{(1-p)^{|\mathcal{M}|}}. \end{aligned}$$

Note that the above condition is equivalent to

$$\theta \leq \inf_{0 \leq u \leq \frac{\sum_{i \in \mathcal{M}} \lambda_i}{F}} \frac{|\mathcal{M}|}{\sum_{i \in \mathcal{M}} \lambda_i - uF} D_{|\mathcal{M}|}(u | p) \quad (23)$$

for all  $\mathcal{M} \subset \{1, 2, \dots, N\}$ . For a fixed  $\mathcal{M}$ , the condition (23) is shown in [11] to be equivalent to

$$\begin{aligned} \sum_{i \in \mathcal{M}} \lambda_i &\leq -\frac{|\mathcal{M}|}{\theta} \log \left[ (1-p)^{|\mathcal{M}|} \right. \\ &\quad \left. + (1 - (1-p)^{|\mathcal{M}|}) \exp\left(-\frac{F\theta}{|\mathcal{M}|}\right) \right]. \end{aligned} \quad (24)$$

Thus, we obtain a lower bound on the effective capacity region as

$$\{\vec{\lambda} \mid \text{Inequality (24) holds for all } \mathcal{M} \subset \{1, 2, \dots, N\}\}. \quad (25)$$

*Remark:* Note that (25) reduces to (16) when all  $\lambda_i$  are equal. Thus, we not only reproduce a lower bound on the effective capacity region for the case with identical offered loads (which is the same as the effective capacity region found in [11]), but also solve the more general problem with non-identical offer loads.

## VI. CONCLUSIONS

In this paper we study the problem of characterizing the delay performance of complex wireless scheduling algorithms. We present a new technique for addressing the complexity issue of the calculus of variations problem involved in the sample-path large deviation approach. Our new technique combines sample-path large deviations with Lyapunov stability, which may develop into a powerful approach to study a large class of scheduling algorithms. We also illustrate the potential of such an approach through an example.

## REFERENCES

- [1] L. Tassiulas and A. Ephremides, "Stability Properties of Constrained Queueing Systems and Scheduling Policies for Maximum Throughput in Multihop Radio Networks," *IEEE Transactions on Automatic Control*, vol. 37, no. 12, pp. 1936–1948, December 1992.
- [2] F. P. Kelly, "Effective Bandwidth in Multiclass Queues," *Queueing Systems*, vol. 9, pp. 5–16, 1991.
- [3] A. I. Elwalid and D. Mitra, "Effective Bandwidth of General Markovian Traffic Sources and Admission Control of High Speed Networks," *IEEE/ACM Transactions on Networking*, vol. 1, no. 3, pp. 329–343, June 1993.
- [4] G. Kesidis, J. Walrand, and C.-S. Chang, "Effective Bandwidth for Multiclass Markov Fluid and other ATM Sources," *IEEE/ACM Transactions on Networking*, vol. 1, no. 4, pp. 424–428, Aug. 1993.
- [5] D. D. Botvich and N. G. Duffield, "Large Deviations, the Shape of the Loss Curve, and Economies of Scale in Large Multiplexers," *Queueing Systems*, vol. 20, pp. 293–320, 1995.
- [6] N. G. Duffield and N. O'Connell, "Large deviations and overflow probabilities for the general single server queue, with application," *Math. Proc. Camb. Phil. Soc.*, vol. 118, pp. 363–374, 1995.
- [7] 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.
- [8] D. Wu and R. Negi, "Effective Capacity: A Wireless Link Model for Support of Quality of Service," *IEEE Transactions on Wireless Communications*, vol. 2, no. 4, pp. 630–643, July 2003.
- [9] A. Eryilmaz and R. Srikant, "Scheduling with Quality of Service Constraints over Rayleigh Fading Channels," in *Proceedings of the IEEE Conference on Decision and Control*, 2004.
- [10] S. Shakkottai, "Effective Capacity and QoS for Wireless Scheduling," available at <http://www.ece.utexas.edu/~shakkott/pub.html>, 2004.
- [11] L. Ying, R. Srikant, A. Eryilmaz, and G. E. Dullerud, "A Large Deviations Analysis of Scheduling in Wireless Networks," in *Workshop on Rare Events in Communication Networks, EURANDOM*, February 2005.

- [12] D. Bertsimas, I. C. Paschalidis, and J. N. Tsitsiklis, "Asymptotic Buffer Overflow Probabilities in Multiclass Multiplexers: An Optimal Control Approach," *IEEE Transactions on Automatic Control*, vol. 43, no. 3, pp. 315–335, March 1998.
- [13] S. Shakkottai, "Modes of overflow, effective capacity and qos for wireless scheduling," in *Proceedings of IEEE International Symposium on Information Theory*, Yokohama, Japan, July 2003.
- [14] A. Dembo and O. Zeitouni, *Large Deviations Techniques and Applications*, 2nd ed. New York: Springer-Verlag, 1998.
- [15] H. K. Khalil, *Nonlinear Systems*, 2nd ed. Upper Saddle River, New Jersey: Prentice-Hall, 1996.
- [16] J. G. Dai, "On Positive Harris Recurrence of Multiclass Queuing Networks: A Unified Approach via Fluid Limit Models," *Annals of Applied Probability*, vol. 5, no. 1, pp. 49–77, 1995.
- [17] A. Shwartz and A. Weiss, *Large Deviations for Performance Analysis: Queues, Communications, and Computing*. London: Chapman & Hall, 1995.