



Management of a Shared-Spectrum Network in Wireless Communications

Shining Wu,^a Jiheng Zhang,^b Rachel Q. Zhang^b

^a Department of Logistics and Maritime Studies, Hong Kong Polytechnic University, Clear Water Bay, Kowloon, Hong Kong; ^b Department of Industrial Engineering and Decision Analytics, Hong Kong University of Science and Technology, Clear Water Bay, Kowloon, Hong Kong

Contact: sn.wu@polyu.edu.hk,  <http://orcid.org/0000-0003-3931-1079> (SW); jiheng@ust.hk,  <http://orcid.org/0000-0003-3025-1495> (JZ); rzhang@ust.hk (RQZ)

Received: February 8, 2015

Revised: October 16, 2015; July 12, 2016; October 3, 2016; February 11, 2017; July 28, 2017

Accepted: October 5, 2017

Published Online in Articles in Advance: July 25, 2018

Subject Classifications: queues: applications, priority; communications

Area of Review: Stochastic Models

<https://doi.org/10.1287/opre.2017.1707>

Copyright: © 2018 INFORMS

Abstract. We consider a band of the electromagnetic spectrum with a finite number of identical channels shared by both licensed and unlicensed users. Such a network differs from most many-server, two-class queues in service systems, including call centers, because of the restrictions imposed on the unlicensed users to limit interference to the licensed users. We first approximate the key performance indicators—namely the throughput rate of the system and the delay probability of the licensed users under the asymptotic regime, which requires the analysis of both scaled and unscaled processes simultaneously using the averaging principle. Our analysis reveals a number of distinctive properties of the system. For example, sharing does not affect the level of service provided to the licensed users in an asymptotic sense even when the system is critically loaded. We then study the optimal sharing decisions of the system to maximize the system throughput rate while maintaining the delay probability of the licensed users below a certain level when the system is overloaded. Finally, we extend our study to systems with time-varying arrival rates and propose a diffusion approximation to complement our fluid one.

Funding: This work was supported by the Hong Kong Research Grants Council [Grants 622713, 16500615, 16501015] and by Hong Kong Polytechnic University [Grants 252019/16E, 1-ZE5G].

Supplemental Material: The e-companion is available at <https://doi.org/10.1287/opre.2017.1707>.

Keywords: spectrum management • many-server queues • fluid approximation • averaging principle

1. Introduction

The radio spectrum refers to the range of frequencies suitable for wireless communications in television and radio broadcasting, aviation, public safety, cell phones, and so on. Until recently, spectrum regulatory bodies, including the Federal Communications Commission (FCC) in the United States and the European Telecommunications Standards Institute, have always allocated spectrum bands exclusively to certain service providers whose users are referred to as primary or licensed users, often based on the radio technologies available at the time of allocation. Such static spectrum allocation mitigates interference to essential services, yet it creates underutilization of the allocated spectrum, which can be below 20% even during high-demand periods in certain geographic areas. For instance, during the high-demand period of a political convention held in New York City in 2004, only about 13% of the allocated spectrum was utilized (Prasad et al. 2010). Studies conducted by the FCC, universities, and industry also revealed that a major part of the spectrum is not fully utilized most of the time. On the other hand, over the past decades, the convergence of voice and data in wireless communications triggered by the convergence of wireless and Internet technologies has

led to an explosion in the number of bits transmitted over the air (Biglieri et al. 2013). Since it is usually difficult to open up higher-frequency bands for mobile applications as transmission becomes less reliable in those bands, the existing radio spectrum for data transmission is reaching its capacity.

A natural approach to alleviate the artificial scarcity of spectrum resulting from static allocation is to allow opportunistic use of temporarily idle channels by unlicensed or secondary users to increase the throughput of already-allocated spectrum. This is referred to as opportunistic spectrum access (Hossain et al. 2009). However, allowing unlicensed users access may cause interference to existing licensed users. Thus, such a paradigm of operation requires (1) the knowledge of the state of frequency bands (e.g., channel availability, queues) in real time and (2) an effective control mechanism to govern spectrum usage by unlicensed users, which has led to the development of the concept of cognitive radio, first introduced by Mitola and Maguire (1999). Using advanced radio- and signal-processing technology, cognitive radio is a software-defined radio device that can intelligently sense and explore the spectrum environment, track changes, communicate information among different transceivers, and react

according to a control mechanism (Hossain et al. 2009). It is widely regarded as one of the most promising technologies for future wireless communications and may potentially mitigate, through dynamic spectrum access, the problem of radio spectrum scarcity.

It is obvious that implementation of a cognitive radio network involves both technological and operational issues, yet much of the research is focused on the former (see Section 2.1 for some relevant literature). In this paper, we focus on the operational issues by considering a band of spectrum with multiple identical channels shared by both licensed and unlicensed users. Since the spectrum has already been allocated to the licensed users, and it is usually difficult to set aside a subset of channels for either group in reality for technical reasons, we assume all the channels are accessible by both licensed and unlicensed users as in most existing literature in electrical engineering. Furthermore, although concurrent transmission is allowed in some networks under which the main concern is technological (e.g., the power level at which an unlicensed user is allowed to transmit), we focus on systems where each channel serves only one user at a time, referred to as the interweave paradigm (Biglieri et al. 2013). Thus, the network considered is a two-class queue served by a single pool of homogeneous servers as in applications in service systems, such as call centers and healthcare, but with some distinctive features as a result of the restrictions imposed on the unlicensed users (Hossain et al. 2009). (1) When all the channels are occupied upon arrival, a licensed user will join a queue along with other waiting licensed users who will be served first in, first out (FIFO) as soon as a channel becomes available, while an unlicensed user will join a queue along with other waiting unlicensed users and will only be allowed to sense channel availability *periodically*. An unlicensed user can occupy a channel only when an available channel is detected and no licensed users are waiting and may also abandon the system every time he senses but finds no available channel. Such a queue where users wait for retrial is referred to as an orbit queue in the queueing literature and is common in computer and communications networks (Artalejo and Gómez-Corral 2008). (2) When in transmission, a licensed user can transmit until his service requirement is fulfilled, while an unlicensed user is only allocated a fixed amount of time, referred to as a *service session*, approaching the end of which he has to stop transmission to sense the environment as sensing cannot occur simultaneously with data transmission. The unlicensed user will be allowed to continue for another service session only if he senses no waiting licensed users. Otherwise, he has to release the channel and join the orbit queue along with other unlicensed users or abandon the system if he needs more time. Note that data transmission can be interrupted

and resumed; hence more complicated control policies than those in call centers are allowed, which leads to new managerial insights.

Assuming that perfect sensing can be achieved in a fixed amount of time and both licensed and unlicensed users arrive according to Poisson processes, we first perform in-depth analysis on the key performance indicators in the management of shared spectrum networks—namely, the delay probability of the licensed users and the system throughput rate. We then focus on the restrictions that need to be imposed on the unlicensed users when in service and waiting (i.e., the length of a service session and the sensing frequency while waiting). Intuitively, the longer a service session is, the less sensing an unlicensed user needs to perform, and hence a higher system throughput rate. Yet longer service sessions can cause more interference to the licensed users. Likewise, the more frequently an unlicensed user senses channel availability while waiting, the sooner the user is able to find an available channel but the more interference the user causes to the licensed users. Thus, there is a trade-off between the throughput rate and the level of interference to the licensed users when deciding on the length of a service session and the sensing frequency. The goals of this research are to answer the following questions: (1) Should a given band of spectrum be shared with unlicensed users? (2) When sharing is permitted, how long should unlicensed users be allowed to transmit each time they occupy a channel, and how frequently should they be allowed to sense channel availability while waiting? (3) Under what conditions is sharing more beneficial? (4) How will the decision change with uncertain arrival rates or time-varying arrivals?

Since the band of spectrum considered usually consists of hundreds or thousands of channels, we can treat the system as a large network and approximate the performance under the asymptotic regime as in Gupta and Kumar (2000) and El Gamal et al. (2006). Because of the restrictions imposed on the unlicensed users when in service and waiting, we need to analyze both scaled and unscaled processes *simultaneously* using the averaging principle (i.e., approximating the unscaled process by its long-run average). We then formulate the problem as finding the optimal restrictions on the unlicensed users to maximize the throughput rate while maintaining the delay probability of licensed users below a certain level. Our main findings are as follows:

1. *Sensing frequency of the unlicensed users while waiting*: Surprisingly, sensing frequency does not affect the system performance asymptotically as long as the unlicensed users are required to sense channel availability, which takes time and prevents them from occupying idle channels instantaneously. Thus, there is no need to impose any restriction on the sensing frequency

from the operational perspective. The decision thus should primarily be based on technological concerns—for instance, power consumption associated with each sensing activity.

2. *The length of a service session:* Intuitively, shorter service sessions should cause less interference to and hence lower the delay probability of the licensed users. However, with shorter service sessions, the unlicensed users need more service sessions to finish their service and hence need to perform more sensing activities while occupying a channel. Thus, shorter service sessions do not always improve the delay probability.

3. *Optimal sharing decisions:* When the system is underloaded or critically loaded, the interference of the unlicensed users to the licensed users is negligible, and there is no need to impose a restriction on the service process of the unlicensed users either. That is, allowing the unlicensed users to complete their transmissions without restriction will not cause any interference to licensed users asymptotically as the delay probability is zero. This result is very different from that of most non-preemptive queueing systems under which the delay probability is strictly between 0 and 1 when the system is critically loaded.

When the system is overloaded, the delay probability of the licensed users is quasi-convex in the length of the service sessions of the unlicensed users, strictly between 0 and 1 and increasing in the load. Thus, a restriction on the service process of the unlicensed users should be imposed only when the load is above a threshold. Furthermore, a shorter service session should be allocated as the load increases until spectrum sharing is no longer feasible.

The insight that it is possible to improve spectrum utilization while guaranteeing a very high service level, expected by licensed users in practice, is very encouraging news. Thus, spectrum sharing can potentially be a socially optimal solution to alleviating spectrum scarcity.

4. For a given system load, a shorter service session should be allocated to the unlicensed users (1) as the proportion of the licensed users increases, (2) if there are fewer licensed users with longer service times, or (3) if there are more unlicensed users with shorter service times. As the service session shortens, more unlicensed users will abandon the system, which lowers the throughput rate under scenarios (1) and (2). Therefore, spectrum sharing is beneficial to systems with a smaller proportion of licensed users or a large number of licensed users with shorter service times.

5. When the arrival rates are time varying, a shorter service session should be allocated to the unlicensed users during busy periods. Although optimal control requires continuous adjustment in real time, near-optimal control can be accomplished with occasional adjustments.

To the best of our knowledge, this is the first comprehensive study of a shared network in wireless communications. Although there have been some attempts by researchers in electrical engineering using relatively simple queueing models, our model captures many more of the features of such a system. We are able to uncover complicated system dynamics and obtain managerial insights different from those drawn from the many well-studied service systems. Our work not only opens the door for new applications of existing queueing theory in wireless communications but also may stimulate the development of new methodologies.

The remainder of this paper is organized as follows. We review the relevant literature in both electrical engineering and queueing theory in the next section and describe the problem of dynamic spectrum sharing in detail in Section 3. In Section 4, we provide a fluid approximation and study the optimal sharing decisions of the system. In Section 5, we offer the intuition behind the construction of the fluid model and give justifications for the fluid approximation. We extend our analysis to systems with time-varying arrival rates and discuss a diffusion-scaled approximation in Section 6. We conclude our paper and provide some future research directions in Section 7. The proofs can all be found in the online appendix.

2. Literature Review

In this section, we first provide some background on the research on opportunistic spectrum access, mostly in electrical engineering. Since we model a shared network as a multiclass, many-server queue where the unlicensed users join an orbit queue and analyze it using the averaging principle, we review the relevant literature in queueing theory and its applications.

2.1. On Opportunistic Spectrum Access

Most of the work on opportunistic spectrum access focuses on the technological issues such as the sensing technology to detect idle channels (Mishra et al. 2006), signal encoding (Devroye et al. 2006), and the control of the transmit power to limit interference (Bansal et al. 2008). For research on various technological issues associated with cognitive radio, readers may refer to Akyildiz et al. (2006) and Goldsmith et al. (2009).

Research on the operational issues under simplified settings, however, remains scant. Huang et al. (2008) performed an analytical study on a single-channel system with one licensed and one unlicensed user, as well as numerical studies on a multichannel system. They also consider the decisions on the sensing frequency of unlicensed users and how long unlicensed users should be allowed to transmit in their numerical study. Zhao et al. (2008) studied the optimal access strategy of an unlicensed user based on the sensing outcome given that each channel has already been assigned to

a specific licensed user, while Capar et al. (2002) compared the system performance in terms of bandwidth utilization and blocking probability when a licensed user can be assigned to any channel randomly or in a controlled way.

For a more comprehensive picture of the various issues in dynamic spectrum management and cognitive radio networks, readers may refer to Hossain et al. (2009) and Biglieri et al. (2013).

2.2. On Queueing Theory and Applications

2.2.1. Multiclass, Many-Server Queues. Since a band of spectrum consists of hundreds or thousands of channels and there are both licensed and unlicensed users, the literature of multiclass, many-server queues is relevant. The study of many-server queues was substantiated by the seminal work of Halfin and Whitt (1981), who derive the steady-state distribution of the diffusion limits and establish the square root law describing the relationship between the system load and delay probability. The mathematical insights of the square root law have since been extended and widely adopted in the daily management of call centers around the world. Later, Puhalskii and Reiman (2000) extended the study to multiclass models.

There is a large body of work on multiclass, many-server systems because of their applications in call centers, manufacturing, and computer-communication systems with a focus on asymptotic optimal control of the underlying systems. For example, Atar et al. (2004) studied asymptotic optimal schedule policies; Gurvich and Whitt (2009) proposed a family of queue-and-idleness-ratio rules for routing and scheduling; and Maglaras and Zeevi (2004, 2005) examined the pricing, capacity sizing, and admission control decisions in a differentiated service system with guaranteed (high-priority) and best-effort (low-priority) users. Our model differs from the existing work in that the service (i.e., data transmission) of the unlicensed users may be fulfilled after multiple interruptions, which is not the case in most other applications.

Since the service of unlicensed users may be interrupted by waiting licensed users, the literature on queues with service interruption caused by preemptive priority, which dates back to White and Christie (1958) in single-server settings, is also relevant. For a review on some of the early work, we refer the reader to Jaiswal (1968). Among the existing work, most focuses on characterizing the steady-state distributions of the queue length, the sojourn time, and so on for a given priority discipline. For example, Brosh (1969) derived the expressions for the expected time from arrival to inception of service and provides bounds for the expected sojourn time for each class when all classes have the same service rates. Buzen and Bondi (1983) obtained the exact expressions for the mean sojourn times when

all classes have the same service rates and provide approximations when different classes have different service rates. Recently, Wang et al. (2015) conducted the exact analysis of the steady state of a preemptive $M/M/c$ queue when different classes have different service rates. In our paper, we focus on the control of the service process of the unlicensed users—that is, how their service processes should be interrupted.

2.2.2. Orbit Queues. Since the unlicensed users join an orbit queue in our setting, the literature along this line is also relevant. Yang and Templeton (1987) and Falin and Templeton (1997) offered a survey and a comprehensive summary of the earlier papers, respectively. Later, Mandelbaum et al. (2002) provided an analytical approximation to the key performance of a many-server queueing system with abandonment and retrials under an asymptotic regime. In all these papers, even though customers may join an orbit queue for retrial if they cannot be served immediately upon arrival, their service cannot be interrupted once started.

Recently, a number of studies considered systems where customers may require repeat service as a result of unresolved or new issues. For instance, de Véricourt and Zhou (2005) and Zhan and Ward (2014) studied a customer-routing problem in call centers with callbacks, while de Véricourt and Jennings (2008) and Yom-Tov and Mandelbaum (2014) examined a staffing problem for membership services and healthcare systems where customers may require multiple rounds of service. These systems differ from ours in that customers will wait in a FIFO queue for retrial if the systems are busy upon arrival, although they will first join an orbit queue after they have had a round of service. Allowing the unlicensed users to retry and join an orbit queue as in our setting significantly complicates the analysis since there may be a large number of customers switching frequently between being in service and being in the orbit queue.

2.2.3. The Averaging Principle. Only a few studies in the queueing literature have required the use of the averaging principle. Building on a fundamental theory of the averaging principle by Kurtz (1992), Hunt and Kurtz (1994) studied martingales and related random measures of large loss networks. Whitt (2002) summarized the early studies on scheduling multiclass queues using the averaging principle. Recently, a series of studies by Perry and Whitt (2011a, b, 2013) applied the averaging principle to obtain both the fluid and diffusion limits for an overloaded X model of many-server queues and to derive insights about the asymptotic optimal control of the system. Pang and Perry (2015) applied the averaging principle to obtain a logarithmic safety staffing rule for call centers with call blending. We adopt some of the methodologies developed by Hunt and Kurtz (1994) and Perry and Whitt (2011a).

3. Problem Description and Assumptions

3.1. The Sharing Network and Performance Measures

We consider a band of spectrum consisting of n identical channels shared by both licensed and unlicensed users, denoted as user types 1 and 2, respectively, and each channel can only be occupied by one user at a time. That is, concurrent transmission is not allowed. Furthermore, we assume that perfect sensing can be achieved in a fixed amount of time $1/\mu_s$ by an unlicensed user. Type i users arrive according to a Poisson process with the rate λ_i^n and require an exponential amount of service time with the rate μ_i , $i = 1, 2$. If there is an available channel, an arriving licensed user will occupy it immediately until his service requirement is fulfilled. Otherwise, the user will join a queue along with other waiting licensed users who will be served FIFO as the channels become available.

Next, we describe the service and waiting processes of the unlicensed users in the shared network illustrated in Figure 1, where $I^n(t)$ is the number of idle channels, and $Q_i^n(t)$ is the queue length of type i users at time t . According to the policy,

$$Q_1^n(t)I^n(t) = 0. \quad (1)$$

Upon arrival, an unlicensed user will occupy a channel if there is one available. Otherwise, he will join an orbit queue along with other waiting unlicensed users with probability $1 - \phi$ or abandon the system.

- *The service process:* Once he occupies a channel, an unlicensed user is allocated a fixed amount of uninterrupted time, referred to as a *service session* (Liu and

Wang 2010), regardless of his service requirement. If he needs more time and finds no licensed user waiting at the end of a session through sensing, he is allowed to continue for another service session. Since sensing cannot occur simultaneously with data transmission and must be interweaved, he needs to devote the last $1/\mu_s$ amount of time in each service session to sense the environment if he needs more time. Hence, we denote the length of a service session by $1/\mu_t + 1/\mu_s$, where $1/\mu_t$ is the amount of time allowed for transmission in a service session. If the unlicensed user completes his transmission within $1/\mu_t$ amount of time in a session, he will release the channel without sensing and leave the system. Otherwise, the user will have to sense the environment, and his service will be interrupted if he finds a waiting licensed user, in which case he will join the orbit queue with probability $1 - \phi$ or abandon the system.

- *The waiting process:* While waiting in the orbit queue, an unlicensed user will only be allowed to sense channel availability periodically. Let $1/\theta$ denote the time between sensing activities, which includes the time needed for sensing channel availability. After each sensing activity, the user will occupy a channel if he finds an idle one. Otherwise, the user will abandon the system with probability ϕ or stay in the orbit queue for another sensing activity with probability $1 - \phi$.

As one can see, the network has its distinctive characteristics, which are not present in most existing multiclass, many-server queueing systems because of the restrictions on the service and waiting processes of the unlicensed users—that is, the transmission time $1/\mu_t$

Figure 1. (Color online) The Spectrum Sharing Network

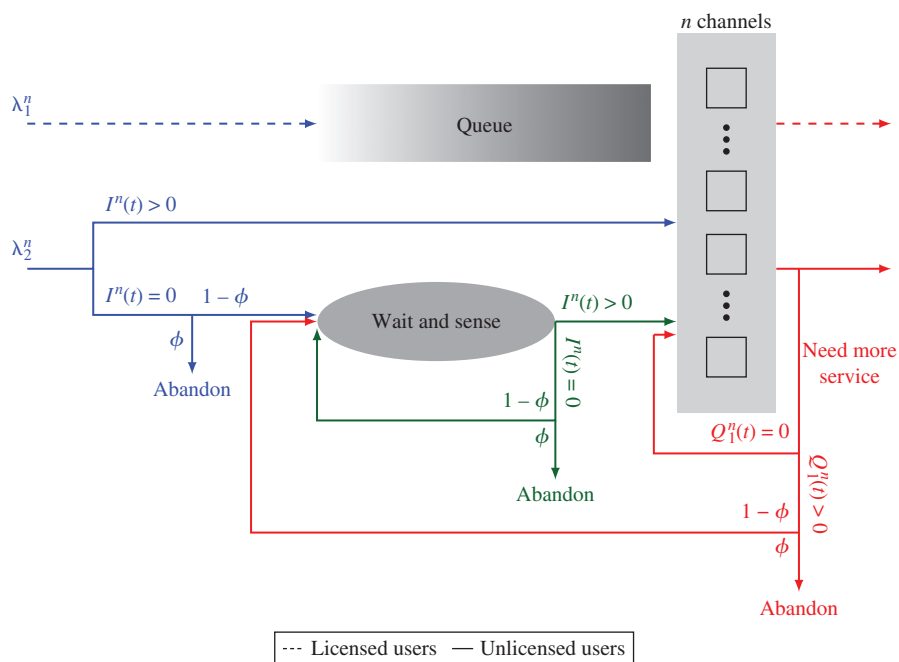


Table 1. Comparison of Performance Measures with Deterministic vs. Exponential Times

$1/\mu_t$	n	Delay probability		Throughput rate	
		Deterministic	Exponential	Deterministic	Exponential
∞	100	0.2528 ± 0.0036	0.2531 ± 0.0044	0.7682 ± 0.0014	0.7678 ± 0.0024
	500	0.2120 ± 0.0014	0.2142 ± 0.0019	0.7928 ± 0.0007	0.7915 ± 0.0012
	1,000	0.2075 ± 0.0017	0.2075 ± 0.0018	0.7957 ± 0.0008	0.7957 ± 0.0009
	2,000	0.2045 ± 0.0009	0.2044 ± 0.0008	0.7975 ± 0.0005	0.7974 ± 0.0005
	4,000	0.2014 ± 0.0010	0.2021 ± 0.0008	0.7992 ± 0.0004	0.7988 ± 0.0004
	Fluid			0.1995	0.8000
0.6	100	0.2360 ± 0.0040	0.2314 ± 0.0036	0.7656 ± 0.0021	0.7662 ± 0.0022
	500	0.1979 ± 0.0009	0.1953 ± 0.0011	0.7901 ± 0.0007	0.7899 ± 0.0005
	1,000	0.1906 ± 0.0013	0.1892 ± 0.0007	0.7949 ± 0.0006	0.7939 ± 0.0005
	2,000	0.1872 ± 0.0006	0.1854 ± 0.0010	0.7968 ± 0.0004	0.7964 ± 0.0005
	4,000	0.1854 ± 0.0003	0.1838 ± 0.0006	0.7979 ± 0.0002	0.7973 ± 0.0003
	Fluid		0.1813		0.7987
0.2	100	0.2262 ± 0.0031	0.2259 ± 0.0033	0.7643 ± 0.0023	0.7641 ± 0.0019
	500	0.1916 ± 0.0024	0.1918 ± 0.0020	0.7878 ± 0.0014	0.7871 ± 0.0011
	1,000	0.1860 ± 0.0010	0.1855 ± 0.0010	0.7914 ± 0.0005	0.7914 ± 0.0007
	2,000	0.1820 ± 0.0011	0.1820 ± 0.0010	0.7940 ± 0.0007	0.7938 ± 0.0006
	4,000	0.1802 ± 0.0005	0.1804 ± 0.0007	0.7954 ± 0.0003	0.7949 ± 0.0004
	Fluid		0.1784		0.7960

in a service session and the sensing frequency θ . The data transmission of the unlicensed users can be interrupted and resumed for any number of times, and sensing for channel availability by the unlicensed users in the queue is only allowed periodically. As a result, an unlicensed user may abandon the system upon arrival, after spending some time in the queue without receiving any service, or after receiving partial service. Furthermore, each unlicensed user in the orbit queue needs to sense channel availability independently, which guarantees certain idleness in the system even when there are waiting unlicensed users. These features are new in the queuing literature and interesting, yet they significantly complicate the analysis.

The performance measures we are concerned with are the throughput rate of the system and the probability that all the channels are occupied upon the arrival of a licensed user, referred to as the delay probability. The goal is to find the transmission time $1/\mu_t$ in a service session and the sensing frequency θ of the unlicensed users that maximize the throughput of the unlicensed users while guaranteeing the delay probability of the licensed users below a certain level.

3.2. Modeling Assumptions

Since the problem is analytically intractable, we first approximate the deterministic transmission time, sensing time, and the time between consecutive sensing activities in the orbit queue by the exponential distributions with the same means. Table 1 presents a simulation study of the delay probability of the licensed users and the throughput rate of the unlicensed users with deterministic times and exponential times when $\lambda_1 = 0.2$, $\lambda_2 = 0.9$, $1/\mu_1 = 1/\mu_2 = 1$, $1/\mu_s = 0.001$, $\theta = 0.4$, and $\phi = 0.5$. For $1/\mu_t \in \{\infty, 0.6, 0.2\}$, we set $n = 100, 500, 1,000, 2,000, 4,000$, and let $\lambda_i^n = n\lambda_i$. We report the

means and 0.95 confidence intervals of the delay probabilities and throughput rates. As one can see, approximating the deterministic times by the exponential times does not reduce the accuracy very much, especially when n is large as in our application, where n is in the hundreds or thousands.

With the exponential times mentioned above, the probability that an unlicensed user will complete his transmission in a service session is given by $p = \mu_2/(\mu_2 + \mu_t)$. Furthermore, the actual amount of time an unlicensed user will occupy a channel in each service session follows a phase-type distribution with mean

$$\frac{1}{\mu} = \frac{1}{\mu_2 + \mu_t} + (1-p) \cdot \frac{1}{\mu_s} = \frac{\mu_t + \mu_s}{(\mu_2 + \mu_t)\mu_s}, \quad (2)$$

which is less than the allocated session time $1/\mu_t + 1/\mu_s$. If we let $Z_i^n(t)$ denote the number of channels occupied by type i users at time t , the instantaneous throughput rate at time t is given by $p\mu Z_2^n(t)$.

With hundreds or thousands of channels in a band of spectrum, performing an analytical study of the shared network under a large system scaling to be defined below in Definition 1 is not only for technical tractability but also appropriate.

Definition 1 (Asymptotic Regime). There exist positive real numbers λ_i , $i = 1, 2$, such that

$$\lim_{n \rightarrow \infty} \frac{\lambda_i^n}{n} = \lambda_i \quad \text{and} \quad \frac{\lambda_1}{\mu_1} < 1.$$

Here, λ_i represents the size of type i users. Different cognitive radio networks have different proportions of licensed and unlicensed users. In IEEE 802.22 wireless regional area networks, unlicensed users outnumber licensed users (Zhang et al. 2009, Jia et al. 2008) (i.e., $\lambda_2 > \lambda_1$), while in TV white space networks, licensed

users are the majority (van de Beek et al. 2012) (i.e., $\lambda_1 > \lambda_2$). In Gong et al. (2015), the licensed users (from a down-link cellular system) and the unlicensed users (from an ad hoc network) have comparable numbers (i.e., $\lambda_1 \approx \lambda_2$).

Under the asymptotic regime, we add a bar to the existing notation to represent the scaled processes in our model (e.g., $\bar{Q}_i^n(t) = Q_i^n(t)/n$) and use the lower-case form (e.g., $q_i(t)$) to represent the corresponding fluid model, which is proven to be the fluid limit of the scaled processes.

4. Main Results and Insights

Under the asymptotic regime, the processes involved are scaled and then approximated by tractable ones that preserve the relevant information about the system performance. As in most multiclass queueing systems, the queue length of the licensed users, who have a higher priority, will vanish asymptotically. This is not a problem if the queue length of the licensed users does not affect the users in service in an asymptotic sense, which is the case in most applications, and one can still obtain the managerial insights by analyzing the limit of scaled processes alone. However, whether the number of waiting license users is asymptotically small or exactly zero is important in our setting, as it determines whether an unlicensed user should vacate a channel, but the scaled processes fail to preserve such important information. Thus, the analysis requires information from both scaled and unscaled processes, involves tracking the two processes simultaneously, and needs to use the averaging principle. These requirements are rare in the literature with only a few exceptions, such as Perry and Whitt (2011a), Luo and Zhang (2013), and Pang and Perry (2015).

In this section, we first introduce our fluid model $x(t) = (z_1(t), q_1(t), z_2(t), q_2(t))$, which is used to approximate the stochastic process $X^n(t) = (Z_1^n(t), Q_1^n(t), Z_2^n(t), Q_2^n(t))$ in our system with the justifications to be provided in Section 5. We then derive the steady-state performance and study the optimal sharing decisions of the system in the steady state using the fluid approximations.

4.1. The Fluid Model

Definition 2 (Fluid Model). The process $x(t) = (z_1(t), q_1(t), z_2(t), q_2(t))$ evolves according to the constraint

$$0 = [1 - z_1(t) - z_2(t)]q_1(t) \quad (3)$$

and the following differential equations:

$$z_1'(t) = [1 - \beta(t)]\lambda_1 + \alpha(t)[\mu_1 z_1(t) + \mu z_2(t)] - \mu_1 z_1(t), \quad (4)$$

$$q_1'(t) = \beta(t)\lambda_1 - \alpha(t)[\mu_1 z_1(t) + \mu z_2(t)], \quad (5)$$

$$z_2'(t) = [1 - \beta(t)][\lambda_2 + \theta q_2(t)] - [p + \alpha(t)(1 - p)]\mu z_2(t), \quad (6)$$

$$q_2'(t) = (1 - \phi)\beta(t)[\lambda_2 + \theta q_2(t)] + (1 - \phi)\alpha(t)(1 - p)\mu z_2(t) - \theta q_2(t), \quad (7)$$

where $\beta(t)$ and $\alpha(t)$ depend on how constraint (3) is met. If $q_1(t) > 0$, then $\beta(t) = \alpha(t) = 1$; if $z_1(t) + z_2(t) < 1$, then $\beta(t) = \alpha(t) = 0$. Otherwise, setting $A = \lambda_1 + \lambda_2 + \theta q_2(t)$, $B = \mu_1 z_1(t) + \mu z_2(t)$, and $C = \mu_1 z_1(t) + p\mu z_2(t)$, we have

$$\beta(t) = \min\left\{\left(\frac{(A - C)B}{AB - \lambda_1 C}\right)^+, 1\right\}, \quad (8)$$

$$\alpha(t) = \min\left\{\frac{\lambda_1 \beta(t)}{\mu_1 z_1(t) + \mu z_2(t)}, 1\right\}. \quad (9)$$

The fluid model defined above is built on the evolution of the system described in Section 3. As we explain in Section 5 and define formally in Online Appendix B, $\beta(t)$ is the instantaneous delay probability of the licensed users, and $\alpha(t)$ is the instantaneous probability that an unlicensed user has to release the channel after a service session (i.e., there are waiting licensed users in the system), referred to as the interruption probability, under the fluid model. Thus, the differential Equations (4)–(7) are quite intuitive. Take Equation (4), for example. The rate of increase in $z_1(t)$ consists of two parts: (1) When the licensed users arrive (at the rate λ_1), there is an available channel (with probability $1 - \beta(t)$). (2) When the licensed users finish service (at the rate $\mu_1 z_1(t)$) or the unlicensed users finish a service session (at the rate $\mu z_2(t)$), there are waiting licensed users (with probability $\alpha(t)$). The rate of decrease in $z_1(t)$ is $\mu_1 z_1(t)$, which is the rate the licensed users occupying the channels finish service. For Equation (7), the rate of increase in $q_2(t)$ consists of two parts: (1) When the unlicensed users arrive or those in the orbit queue perform sensing (at the rate $\lambda_2 + \theta q_2(t)$), they find all channels occupied (with probability $\beta(t)$) but decide not to abandon the system (with probability $1 - \phi$). (2) When the unlicensed users finish a service session (at the rate $\mu z_2(t)$), they need another one (with probability $1 - p$) and find licensed users waiting (with probability $\alpha(t)$) but do not abandon the system (with probability $1 - \phi$). The rate of decrease in $q_2(t)$ is $\theta q_2(t)$, which is the rate the unlicensed users in queue sense for available channels.

When $z_1(t) + z_2(t) < 1$ or $q_1(t) > 0$, the system dynamics are quite simple and resemble that of the many-server queues in call center applications. For example, when $z_1(t) + z_2(t) < 1$, the differential equations (4)–(7) reduce to $q_1(t) \equiv 0$, and

$$\begin{aligned} z_1'(t) &= \lambda_1 - \mu_1 z_1(t), \\ z_2'(t) &= \lambda_2 + \theta q_2(t) - p\mu z_2(t), \\ q_2'(t) &= -\theta q_2(t). \end{aligned}$$

Otherwise, the system dynamics are more complicated. Moreover, the process $x(\cdot)$ can move back

and forth among different cases, which makes the analysis even more challenging as shown in Online Appendix A.

Despite the complexity, the fluid model can be solved numerically. Furthermore, we can obtain the steady state of the fluid model in Theorem 1 to approximate the steady state of the original system. For example, $\beta := \lim_{t \rightarrow \infty} \beta(t)$ and $TH_2 := \lim_{t \rightarrow \infty} p\mu z_2(t)$ can be used to accurately approximate the steady-state delay probability of the licensed users and the throughput rate of the unlicensed users, respectively. Note that fluid models fail to yield probabilistic performance measures in most applications. Similar to Gurvich and Perry (2012), our fluid model actually provides accurate approximations for them.

4.2. The Steady State of the Fluid Model

While the *offered load* of such a system is $\lambda_1/\mu_1 + \lambda_2/\mu_2$, the *effective load* is endogenous, as the average time for which an unlicensed user occupies a channel $1/\mu$ defined in (2) depends on the decision $1/\mu_t$. Since $1/p$ is the average number of service sessions needed to fulfill the service requirement of an unlicensed user, the effective service time of an unlicensed user is $1/(p\mu)$. Thus, the effective load of the system is

$$\frac{\lambda_1}{\mu_1} + \frac{\lambda_2}{p\mu},$$

where $p\mu = \mu_2\mu_s/(\mu_t + \mu_s)$. Note that the effective load is always no less than the offered load and equals the offered load if and only if there is no restriction on the service process of the unlicensed users (i.e., $1/\mu_t = \infty$). The shorter the transmission time in a service session, the more service sessions (and hence sensing) are needed for the unlicensed users to complete their transmissions, and the more congested the system is. Depending on the effective load of the system, the steady states of the fluid limits are given in the next theorem, whose proof can be found in Online Appendix A.

Theorem 1. *There exists a unique solution to the fluid model. (Note that a vector-valued function $x(t)$ is called a solution of the fluid model if it is absolutely continuous on every closed time interval and satisfies Equations (4)–(7) almost everywhere.) Moreover, the limiting behavior of the fluid model as $t \rightarrow \infty$ can be characterized as follows:*

1. *If $\lambda_1/\mu_1 + \lambda_2/(p\mu) > 1$, then we have $\lim_{t \rightarrow \infty} x(t) = (\lambda_1/\mu_1, 0, 1 - \lambda_1/\mu_1, ((1 - \phi)/(\theta\phi))[\lambda_2 - p\mu(1 - \lambda_1/\mu_1)])$, $TH_2 = p\mu(1 - \lambda_1/\mu_1)$, and (β, α) is the unique solution to*

$$\alpha = \frac{\lambda_1}{\lambda_1 + \mu(1 - \lambda_1/\mu_1)}\beta, \quad (10)$$

$$\gamma = \beta + (1 - \beta) \frac{(1 - p)\alpha}{p + (1 - p)\alpha}, \quad (11)$$

$$\lambda_2 \mathbb{E}[\gamma^K] = \lambda_2 - p\mu \left(1 - \frac{\lambda_1}{\mu_1}\right), \quad (12)$$

where $K \geq 1$ follows a geometric distribution with parameter ϕ .

2. *If $\lambda_1/\mu_1 + \lambda_2/(p\mu) \leq 1$, then we have $\lim_{t \rightarrow \infty} x(t) = (\lambda_1/\mu_1, 0, \lambda_2/(p\mu), 0)$, $TH_2 = \lambda_2$, and $\alpha = \beta = 0$.*

We first describe the intuition behind the delay probability β in Equations (10)–(12) before discussing the steady-state behavior in more detail in the next section. Equation (10) is obtained by plugging $\lim_{t \rightarrow \infty} x(t)$ into (9). Note that $1 - \beta$ is also the probability that an unlicensed user will be served upon arrival or after each sensing activity while waiting in the orbit queue, and $(1 - p)\alpha/(p + (1 - p)\alpha)$ is the probability that an unlicensed user in service will be interrupted. Thus, γ in (11) is the probability that an unlicensed user will experience blockage or interruption and hence needs to decide whether or not to abandon the system at least once. Since $K \geq 1$ represents the number of times an unlicensed user needs to decide whether to abandon the system, $\mathbb{E}[\gamma^K]$ is the probability that an unlicensed user will abandon the system. So the left-hand side of (12) can be understood as the abandonment rate of the unlicensed users, while the right-hand side is also the abandonment rate but calculated by subtracting the rate $p\mu(1 - \lambda_1/\mu_1)$ at which unlicensed users complete their service from the total arrival rate λ_2 . Given that $\mathbb{E}[\gamma^K] = \gamma\phi/(1 - \gamma(1 - \phi))$, we actually have a closed-form expression (see (29) in Online Appendix A) for the delay probability β from solving (10)–(12).

Table 1 also presents a comparison between the simulated delay probability and throughput rate and the approximation based on the fluid model. As one can see, the fluid approximation works well, especially when n is large, which is the case in our application. Furthermore, our simulation also reveals that the average queue length of the licensed users is indeed quite short (vanishes asymptotically). For instance, the 0.95 confidence interval of the queue length of the licensed users is 0.0172 ± 0.0001 with deterministic times and 0.0211 ± 0.0001 with exponential times when $n = 4,000$ and $1/\mu_t = 0.6$.

From Theorem 1, we can see that the system performance is insensitive to θ , the frequency at which the unlicensed users sense for an available channel while waiting in the orbit queue. This is because, although θ affects the transient of the differential equations (4)–(7), it influences the steady state through the total sensing speed θq_2 (i.e., when the derivatives of the left hand side equal 0). As θ increases, the unlicensed users are allowed to sense channel availability more frequently and hence abandon the system sooner, which lowers the number of waiting unlicensed users q_2 . It turns out that, under such a mechanism, the total sensing speed remains constant as θ varies. The insensitivity of θ on the system performance is further confirmed by a simulation study in Online Appendix C. Thus, the decision on the sensing frequency should be based on technological (e.g., power consumption as sensing consumes power) rather than operational concerns.

When the system is effectively underloaded or critically loaded, in which case the offered load is $\lambda_1/\mu_1 + \lambda_2/\mu_2 \leq 1$, all users will be served without delay in the steady state, and no restriction needs to be imposed on the unlicensed users. When the system is effectively overloaded, in which case the offered load may or may not be above 1, only $p\mu(1 - \lambda_1/\mu_1)$ of the unlicensed users will finish service per unit time, and the unlicensed users will experience interference with a positive probability.

Theorem 1 also reveals some interesting steady-state behavior that differs from that of the fluid models in most applications such as call centers:

1. It is well understood in the queueing literature that, if a system is critically loaded, there is a positive probability that delay will occur, even with an extra capacity of $O(\sqrt{\lambda^n})$ in most nonpreemptive models in applications such as call centers (see Halfin and Whitt 1981). In our application, the requirement for unlicensed users to sense channel availability while waiting in the orbit queue guarantees the availability of idle channels for all licensed users upon arrival even when the system is critically loaded, leading to a zero delay probability for licensed users asymptotically. We note a similar result in Pang and Perry (2015) that, by controlling when outbound calls can be made, reserving a logarithmic order number of servers in a call center can achieve a zero delay probability for inbound calls asymptotically when the system is critically loaded.

2. It is also well understood that, when a system is overloaded, customers will experience delay almost surely in most call center applications because all servers are busy all the time (see Whitt 2006). In our application, however, an arriving licensed user still has a chance to enter service upon arrival even when there is a large number of unlicensed users in the orbit queue as it takes time for them to sense channel availability. Hence, the delay probability of the licensed users, which is endogenously determined by the load through (10)–(12), is strictly less than 1. Even if a licensed user is delayed upon arrival, his waiting time is in the order of $O(1/\lambda_1^n)$, which is relatively short but may still be significant in data transmission.

In essence, the restriction that the unlicensed users are not allowed to sense channel availability constantly makes the system operate more like a preemptive one for the licensed users than a nonpreemptive one.

4.3. Sensitivity of the System Performance

By Theorem 1, sensing frequency does not affect the system performance. Thus, we focus on the impact of the length of transmission time $1/\mu_t$ (or, equivalently, the length of the service session) on the throughput of the unlicensed users TH_2 and the delay probability of the licensed users β .

Corollary 1. *Throughput TH_2 is always increasing in the transmission time $1/\mu_t$; that is, allowing the unlicensed users longer service sessions will increase the system throughput rate. The delay probability β is quasi-convex in $1/\mu_t$. More specifically,*

- if $\lambda_1/\mu_1 + \lambda_2/\mu_2 \leq 1$ or $1/\mu_s \geq (1/\mu_2) \cdot (1 - (\mu_2/\lambda_2) \cdot (1 - \lambda_1/\mu_1)) / (1 + (\mu_2/\lambda_1)(1 - \lambda_1/\mu_1))$, then β decreases in $1/\mu_t$ (see Figure 2(a) and 2(b));
- otherwise, there exists a threshold $1/\hat{\mu}_t < \infty$ such that β decreases in $1/\mu_t$ when $1/\mu_t \leq 1/\hat{\mu}_t$ and increases in $1/\mu_t$ when $1/\mu_t > 1/\hat{\mu}_t$ (see Figure 2(c) and 2(d)).

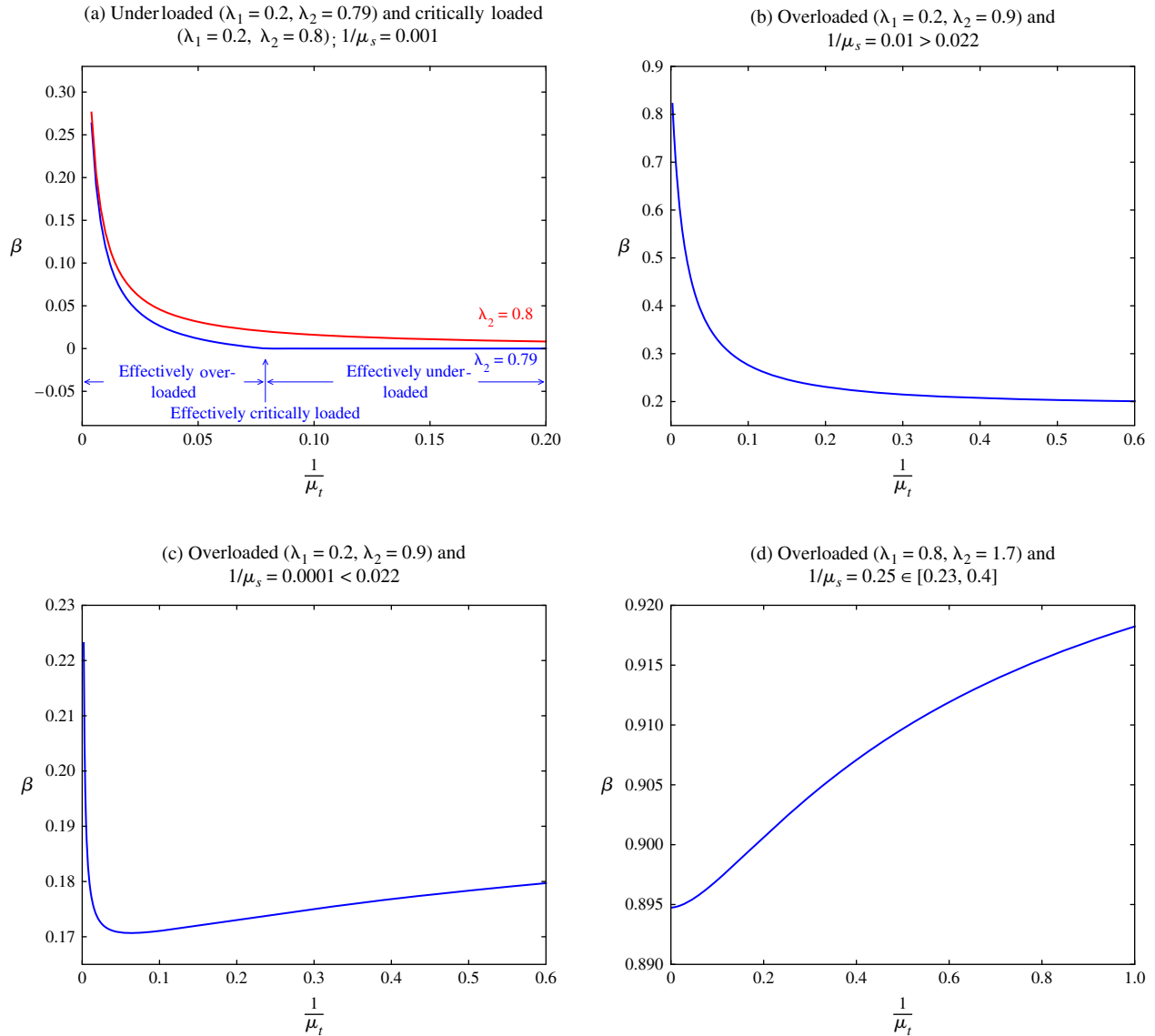
Figure 2 illustrates the delay probability as a function of the transmission time for various λ_1 , λ_2 , and $1/\mu_s$ when $1/\mu_1 = 1/\mu_2 = 1$, $\theta = 0.4$, and $\phi = 0.5$. Note that the purpose of restricting the amount of time the unlicensed users can occupy a channel is to limit the interference of the unlicensed users to the service of the licensed users. Thus, intuitively, shorter service sessions should always lead to a lower delay probability. The corollary reveals that this is true only if $1/\hat{\mu}_t = 0$, which happens when the workload from both types of users are high enough and the sensing time is moderate (see Figure 2(d)). When the system is overloaded and sensing is not too time-consuming, imposing too short service sessions will only increase the effective load and hence the delay probability while imposing relatively longer service sessions will increase the delay probability as expected (see Figure 2(c)). When the system is underloaded or critically loaded, shorter service sessions will either have no impact on the delay probability or turn the system into an effectively overloaded one, increasing the delay probability (see Figure 2(a)). When the system is overloaded and sensing takes a long time, it only makes sense to allow an unlicensed user to transmit for a significant amount of time in order to lower the delay probability (see Figure 2(b)).

4.4. Optimal Sharing Decisions in the Steady State

In this section, we investigate whether a given band of spectrum should be shared with unlicensed users and the transmission time $1/\mu_t$ that maximizes the throughput rate of the unlicensed users while keeping the delay probability of the licensed users below a certain level, η . Note that θ does not affect the system performance by Theorem 1, and the transmission time is the only decision. Furthermore, maximizing the throughput rate of the unlicensed users is equivalent to maximizing the throughput rate of the system since the throughput rate of the licensed users is a constant.

4.4.1. Whether and How to Share. When the system is underloaded or critically loaded, the system may also be effectively overloaded if one allocates shorter service sessions to the unlicensed users. However, by Theorem 1, even if the unlicensed users are allowed to transmit for as long as they need, the delay probability

Figure 2. (Color online) The Delay Probability as a Function of the Transmission Time



converges to zero, and all users are able to complete their transmission without delay as $n \rightarrow \infty$. Thus, we do not need to restrict the service session of the unlicensed users when n is large enough.

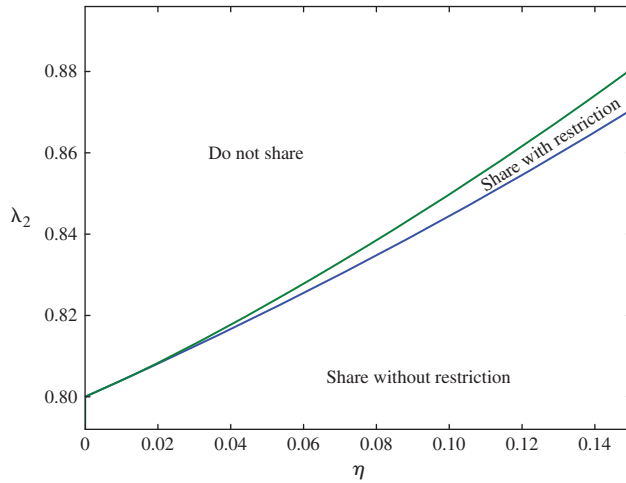
When the system is overloaded, it is also effectively overloaded regardless of the length of the allocated service session. By Theorem 1, $\text{TH}_2 = p\mu(1 - \lambda_1/\mu_1)$ and β is the solution to (12). Thus, the optimization problem can be written as

$$\begin{aligned} \max_{\mu_t \geq 0} \quad & p\mu & (13) \\ \text{s.t.} \quad & \beta \leq \eta, \\ & \mu = \frac{(\mu_2 + \mu_t)\mu_s}{\mu_t + \mu_s}, \\ & p = \frac{\mu_2}{\mu_2 + \mu_t}. \end{aligned}$$

Since the objective function is increasing in $1/\mu_t$, the optimization problem reduces to one of finding the largest $1/\mu_t$ that satisfies the delay constraint. When η is so small that the feasible region is empty, no unlicensed users should be allowed in the system. Once η is large enough to make the feasible region nonempty, unlicensed users will be allowed in the system. As η increases, the optimal transmission time $1/\mu_t^*$ increases. The optimal transmission time $1/\mu_t^* = \infty$ —that is, the unlicensed users are allowed to complete their transmission once they start occupying a channel—if η is larger than the point such that the feasible region becomes unbounded.

Figure 3 demonstrates the optimal spectrum sharing decision as a function of η and λ_2 when $\lambda_1 = 0.2, 1/\mu_1 = 1/\mu_2 = 1, 1/\mu_s = 0.001, \theta = 0.4,$ and $\phi = 0.5$. The upper curve specifies the arrival rate above which the un-

Figure 3. (Color online) The Optimal Sharing Decision as a Function of η and λ_2 for an Overloaded System



licensed users should not be allowed to share the spectrum, and the lower one is the threshold below which there is no need to restrict the service session of the unlicensed users (i.e., $1/\mu_i^* = \infty$).

Since our analysis only holds in an asymptotic sense (as the number of channels n becomes large), there is still a nonnegligible delay probability when the system is underloaded or critically loaded and n is not sufficiently large. For the same example in Figure 3 with $\lambda_1^n/n = 0.2$, Figure 4 demonstrates the optimal sharing decisions, obtained through simulation, as a function of η and λ_2^n/n for $n = 100, 200, 500, 1,000$, in which case the system is underloaded or critically loaded when $\lambda_2^n/n \leq 0.8$. As one can see, the structure of the optimal sharing decision remains the same, and as n increases, sharing is more likely to occur, and the unlicensed users should be allowed longer service sessions.

4.4.2. Sensitivity of the Optimal Decision. The optimal decision $1/\mu_i^*$ and the throughput rate of the unlicensed users TH_2^n depend on the system parameters in the following way.

Proposition 1. *The optimal $1/\mu_i^*$ decreases (i.e., the unlicensed users are allowed a shorter transmission time) as*

- (1) λ_1 increases while keeping $\lambda_1 + \lambda_2$ constant when $\mu_1 = \mu_2$;
- (2) λ_1 and μ_1 decrease while keeping λ_1/μ_1 constant; and
- (3) λ_2 and μ_2 increase while keeping λ_2/μ_2 constant.

Furthermore, the optimal throughput TH_2^n will decrease under (1) and (2).

Note that under all the scenarios, the total offered load $\lambda_1/\mu_1 + \lambda_2/\mu_2$ is kept constant. Proposition 1 states that shorter service sessions should be allocated to the unlicensed users (1) as the proportion of licensed users increases when all users have identical service requirements, (2) if there are fewer licensed users but

with longer service times, and (3) if there are more unlicensed users but with shorter service times. While (1) and (3) are more intuitive, (2) holds because the delay probability only depends on both λ_1/μ_1 and λ_1 . A delay incident of a licensed user is counted as one regardless of his service requirement. With fewer licensed users, each delay contributes more to the delay probability, and it is easy to show that shorter service sessions should be imposed on the unlicensed users.

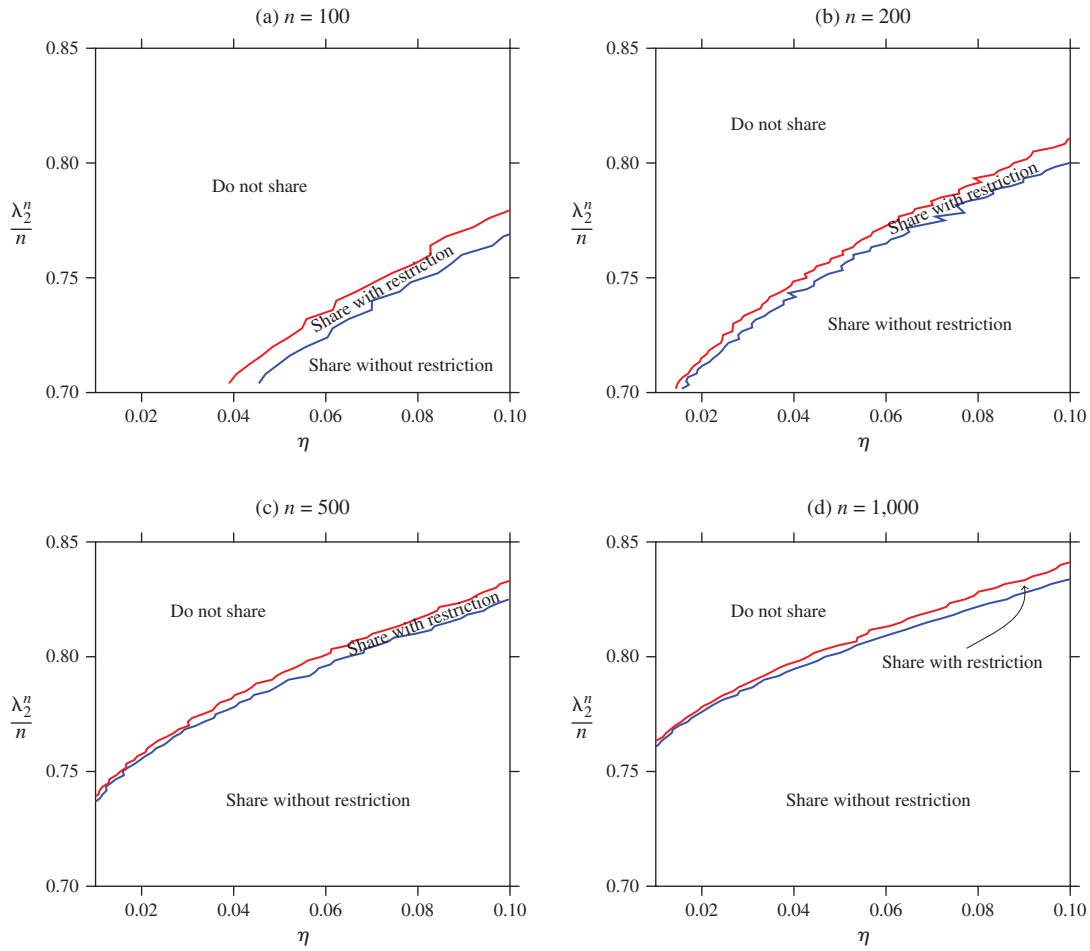
As a result, the optimal throughput rate $p\mu^*(1 - \lambda_1/\mu_1) = (\mu_2\mu_s/(\mu_i^* + \mu_s))(1 - \lambda_1/\mu_1)$ decreases under scenarios (1) and (2) as expected. These suggest that spectrum sharing is beneficial to systems with a smaller proportion of licensed users or a large number of licensed users with shorter service times. Under scenario (3), although shorter service sessions have a negative impact on the throughput rate because of more sensing activities, the increase in the number of unlicensed users with shorter service times has a positive impact. Thus, the impact on throughput rate is not monotone.

5. Justifications for the Fluid Approximation

In this section, we demonstrate in Theorem 2 that the scaled process $\bar{X}^n(t)$ converges to the fluid model $x(t)$ in Section 4. Since the proof of the theorem is quite involved, we describe the main ideas of the proof through the construction of the fluid model, especially the instantaneous delay probability of the licensed users $\beta(t)$ and the instantaneous interruption probability of the unlicensed users $\alpha(t)$. The complete proof can be found in Online Appendix B. For any $T > 0$, let $\mathcal{D}([0, T], \mathbb{R}^4)$ be the space of all right-continuous \mathbb{R}^4 -valued functions on $[0, T]$ with left limits, endowed with the Skorohod J_1 topology. Let “ \Rightarrow ” denote convergence in distribution for random objects in \mathbb{R}^4 equipped with Euclidian topology or $\mathcal{D}([0, T], \mathbb{R}^4)$ with Skorohod J_1 topology.

Theorem 2 (Fluid Approximation). *Under the asymptotic regime, if $\bar{X}^n(0) \Rightarrow x(0)$ as $n \rightarrow \infty$, then $\bar{X}^n(t) \Rightarrow x(t)$ in $\mathcal{D}([0, T], \mathbb{R}^4)$, where $x(t)$ is the fluid model specified in Definition 2.*

Need for Both Scaled and Unscaled Processes. If we let $\Lambda_i^n(t)$ denote the Poisson process with the rate λ_i^n , then $\Lambda_1^n(t + \delta) - \Lambda_1^n(t)$ is the total number of licensed users arriving in a small interval $[t, t + \delta]$, among which $\int_t^{t+\delta} \mathbf{1}_{\{I^n(s)=0\}} d\Lambda_1^n(s)$ will find no idle channels upon arrival and have to wait. Thus, the delay probability of the licensed users during this small time interval is $\mathbb{E}[\int_t^{t+\delta} \mathbf{1}_{\{I^n(s)=0\}} d\Lambda_1^n(s) / (\Lambda_1^n(t + \delta) - \Lambda_1^n(t))]$. That is, the delay probability depends on the information about the unscaled process $I^n(t) \geq 0$ as it determines whether an unlicensed user in service should vacate a channel

Figure 4. (Color online) The Optimal Sharing Decisions as a Function of η and λ_2^n/n 

at the end of a service session. Likewise, we need to keep track of the unscaled process of the queue length of the licensed users $Q_1^n(t)$ and obtain the probability of an unlicensed user in service being interrupted in $[t, t + \delta]$. However, $I^n(t) \geq 0$ vanishes in the asymptotic regime along with the process $Q_1^n(t) \geq 0$ as in most systems with multiple classes, and we need to keep track of both the scaled and unscaled processes in order to obtain the system dynamics and asymptotic system performance.

The System Dynamics Using the Averaging Principle.

To obtain the system dynamics, we need to apply the averaging principle by first expressing the probabilities in $[t, t + \delta]$ as a time average using PASTA (Poisson arrivals see time average). For instance, the fraction of time for which there is no idle channel in the system is

$$\frac{1}{\delta} \int_t^{t+\delta} \mathbf{1}_{\{I^n(s)=0\}} ds = \frac{1}{n\delta} \int_t^{t+n\delta} \mathbf{1}_{\{I^n(t+s-n)=0\}} ds. \quad (14)$$

Let

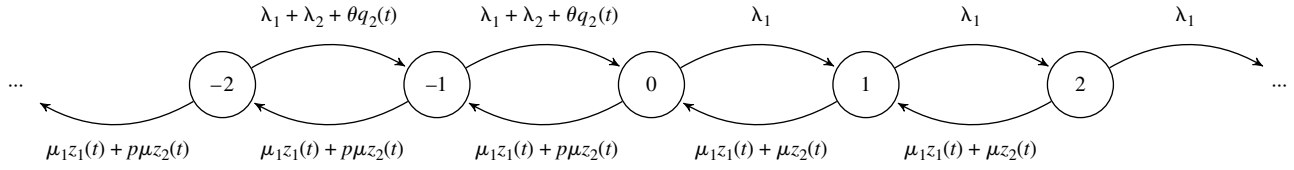
$$m^n(t) = Q_1^n(t) - I^n(t). \quad (15)$$

We study the system dynamics for the unscaled process $m^n(t + s/n)$ for $0 \leq s \leq n\delta$. Note that the process $m^n(t + \cdot/n)$ oscillates around zero on the order of 1. When $m^n(t + s/n) < 0$ (there are idle channels and no licensed users waiting by (1)), the process increases by 1 when there is a new arrival at the rate $\bar{\lambda}_1^n + \bar{\lambda}_2^n$ or one of the unlicensed users in the orbit queue enters service after sensing the system at the rate $\theta \bar{Q}_2^n(t + s/n)$. The process decreases by 1 when a user (licensed or unlicensed) completes service at the rate $\mu_1 \bar{Z}_1^n(t + s/n) + p\mu \bar{Z}_2^n(t + s/n)$. When $m^n(t + s/n) > 0$ (there are licensed users waiting and no idle channels), the process increases by 1 at the rate $\bar{\lambda}_1^n$ and decreases at the rate $\mu_1 \bar{Z}_1^n(t + s/n) + \mu \bar{Z}_2^n(t + s/n)$. We refer readers to Online Appendix B.1 for the detailed system dynamics.

It is the long-run average behavior of $m^n(t + \cdot/n)$ that plays the key role in determining the fraction in (14) when n becomes large in the asymptotic regime.

Explanation for $\beta(t)$ and $\alpha(t)$. As one can see, the process $m^n(t + \cdot/n)$ is not a Markov process since its evolution depends on a higher dimension process than itself. However, if we approximate the abovementioned rates by their fluid counterparts (i.e., $\bar{Z}_i^n(t + s/n)$ by $z_i(t)$,

Figure 5. The Asymptotic Transition Diagram of $m^n(t + \cdot/n)$



$\bar{Q}_2^n(t + s/n)$ by $q_2(t)$, and $\bar{\lambda}_i^n$ by λ_i , we have a Markov process as in Figure 5 whose steady-state distribution π_i can be easily obtained. We use $\pi_i(j)$ to approximate the asymptotic proportion of time for which there are j licensed users in the queue when $j > 0$ and there are $-j$ idle channels when $j < 0$. The delay probability of the licensed users in (14) and the interruption probability of the unlicensed users in the asymptotic regime can be approximated by $\sum_{j=0}^{\infty} \pi_i(j) := \beta(t)$ and $\sum_{j=1}^{\infty} \pi_i(j) := \alpha(t)$, respectively.

6. Extensions

In this section, we extend the problem to systems with time-varying arrival rates and propose a diffusion approximation that can lead to better performance in some cases.

6.1. With Time-Varying Arrival Rates

When the arrival rates vary over time, the optimal decision on the transmission time needs to be adjusted dynamically. Suppose that adjustment of the transmission time can be done instantaneously and the initial state $x(0)$ is given. We can extend the fluid model in Definition 2 to allow time-varying arrivals by adding an argument t to λ_i , μ_t , p , and μ to denote their instantaneous values. Following similar arguments in Online Appendices A and B, we can show that the stochastic processes with time-varying arrival rates converge to the extended fluid model, and there exists a unique solution to time-varying differential equations of the fluid model as long as $\lambda_i(t)$'s are bounded and locally Lipschitz continuous. In this case, the instantaneous throughput rate is $p(t)\mu(t)z_2(t)$, the instantaneous delay probability $\beta(t)$ is given by (8), and the optimization problem over a period of time T can then be written as

$$\begin{aligned} \max_{\mu_i(\cdot)} \quad & \int_0^T p(t)\mu(t)z_2(t) dt \\ \text{s.t.} \quad & \beta(t) \leq \eta, \\ & \mu(t) = \frac{[\mu_2 + \mu_i(t)]\mu_s}{\mu_t(t) + \mu_s}, \\ & p(t) = \frac{\mu_2}{\mu_2 + \mu_t(t)}. \end{aligned}$$

Although such a continuous-time dynamic programming problem can be solved numerically using policy iteration, the resulting policy is hard to implement in practice. Thus, we ask whether a periodically

adjusted policy will work well. As an example, suppose that the arrival rates change over time as in Figure 6: $1/\mu_1 = 1/\mu_2 = 1$ minutes, $1/\mu_s = 0.02$ minute, $\theta = 0.4$, $\phi = 0.5$, and $\eta = 0.2$. Figure 6 plots the transmission times adjusted on an hourly basis and the performance over a 10-hour period. As one can see, our heuristic policy performs very well. According to our numerical experiments, the throughput rate under the hourly adjusted policy is consistently within 0.2% of the optimal throughput rate.

6.2. A Diffusion Scaling

Although our fluid scaling results in good approximations, it leads to a zero delay probability when the system is underloaded and critically loaded, which is not accurate when n is small. Thus, we ask whether a diffusion scaling may work better for underloaded and critically loaded systems.

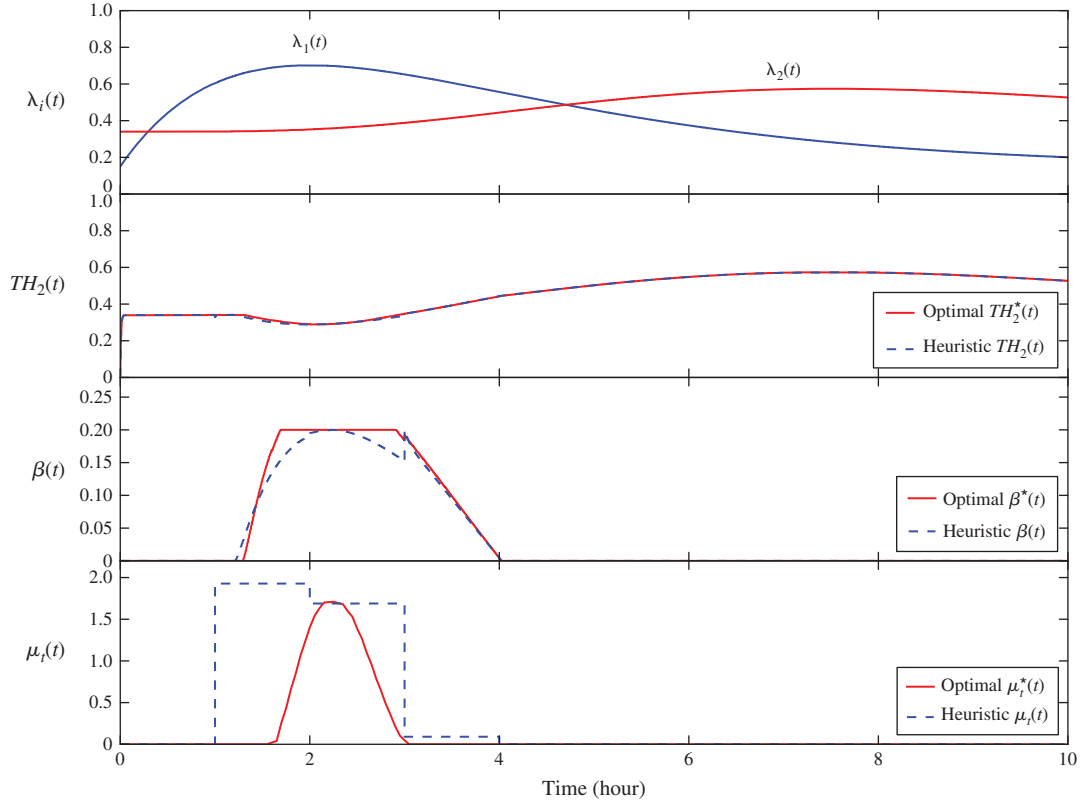
Consider the diffusion scaling where the licensed users grow in the order of $O(\sqrt{n})$ and the unlicensed users in the order of $O(n)$; that is,

$$\begin{aligned} \lambda_1^n &= \tilde{\lambda}_1 \sqrt{n}, \\ \lambda_2^n &= p \mu n + \tilde{\lambda}_2 \sqrt{n}, \end{aligned}$$

and $\tilde{Z}_1^n(t) = Z_1^n(t)/\sqrt{n}$, $\tilde{Z}_2^n(t) = (Z_2^n(t) - n)/\sqrt{n}$, and $\tilde{Q}_i^n(t) = Q_i^n(t)/\sqrt{n}$ are the corresponding diffusion-scaled processes. Such a scaling explicitly assumes that there are far more unlicensed users than licensed users. It can be shown that the diffusion-scaled processes converge, and there exist coefficients C_1, C_2, C_q, C_β , and C_α such that

$$\begin{aligned} \mathbb{E}[Z_1^n(\infty)] &= C_1 \sqrt{n} + o(\sqrt{n}), \\ \mathbb{E}[Z_2^n(\infty)] &= n + C_2 \sqrt{n} + o(\sqrt{n}), \\ \mathbb{E}[Q_1^n(\infty)] &= o(1), \quad \mathbb{E}[Q_2^n(\infty)] = C_q \sqrt{n} + o(\sqrt{n}), \\ a^n &= \frac{C_\alpha}{\sqrt{n}} + o\left(\frac{1}{\sqrt{n}}\right), \quad \beta^n = \frac{C_\beta}{\sqrt{n}} + o\left(\frac{1}{\sqrt{n}}\right). \end{aligned} \quad (16)$$

6.2.1. Estimation of the Coefficients. First, it is easy to see that $C_1 = \tilde{\lambda}_1/\mu_1$, since the licensed users do not abandon. Since the unlicensed users may abandon, the system is stable in the long run, and hence the balance equations are given by letting (4)–(7) be zero and

Figure 6. (Color online) With Time-Varying Arrival Rates and Periodic Adjustments

replacing $(\lambda_i, z_i, q_i, \beta, \alpha)$ by $(\lambda_i^n, \mathbb{E}[Z_i^n(\infty)], \mathbb{E}[Q_i^n(\infty)], \beta^n, \alpha^n)$. Solving the balance equations, we are able to obtain

$$C_\alpha = 0, \quad (17)$$

$$C_\beta = \frac{\theta C_q}{(1-\phi)p\mu}, \quad (18)$$

$$\theta C_q = (1-\phi)[\bar{\lambda}_2 - p\mu C_2 + \theta C_q].$$

It remains to estimate C_2 and C_q . If we are able to derive a closed-form steady-state distribution of the limit of the diffusion-scaled process, we can obtain the value of these coefficients. Although the four-dimensional diffusion-scaled process, $(\tilde{Z}_1^n, \tilde{Q}_1^n, \tilde{Z}_2^n, \tilde{Q}_2^n)$, can be reduced to a three-dimensional process as \tilde{Q}_1^n converges to 0, it has some complicated reflection behavior on the boundary when all channels are busy (i.e., $\tilde{Z}_1^n(t) + \tilde{Z}_2^n(t) = 0$). In general, it is challenging to derive the steady-state distribution of a multidimensional diffusion process, and closed-form expressions of the coefficients are almost impossible. Thus, we propose a heuristic method to derive closed-form approximations for C_2 and C_q and hence C_β .

We pretend that the licensed users occupy $(\bar{\lambda}_1/\mu_1) \cdot \sqrt{n}$ channels exclusively, and the waiting unlicensed users form a steady source of arrival with the rate $\theta C_q \sqrt{n}$. The unlicensed users are served by the remaining $n - (\bar{\lambda}_1/\mu_1)\sqrt{n}$ channels and form an Erlang-B

queue with the arrival rate $\lambda_2^n + \theta C_q \sqrt{n}$ and service rate $p\mu$. In such a network, $\lim_{n \rightarrow \infty} \tilde{Z}_2^n(t) = \tilde{z}_2(t)$ is a reflected Brownian motion with an infinitesimal mean $-p\mu(\tilde{z}_2 - (\bar{\lambda}_2 + \theta C_q)/(p\mu))$ and infinitesimal variance $2p\mu$. Therefore, $\lim_{t \rightarrow \infty} \lim_{n \rightarrow \infty} \tilde{Z}_2^n(t)$ follows a truncated normal distribution with mean $(\bar{\lambda}_2 + \theta C_q)/(p\mu)$ and variance 1 on $(-\infty, -\bar{\lambda}_1/\mu_1)$, and hence

$$C_2 = \mathbb{E}[\lim_{t \rightarrow \infty} \lim_{n \rightarrow \infty} \tilde{Z}_2^n(t)] = \frac{\bar{\lambda}_2 + \theta C_q}{p\mu} - \frac{\Phi(-\bar{\lambda}_1/\mu_1 - (\bar{\lambda}_2 + \theta C_q)/(p\mu))}{\Phi(-\bar{\lambda}_1/\mu_1 - (\bar{\lambda}_2 + \theta C_q)/(p\mu))}, \quad (19)$$

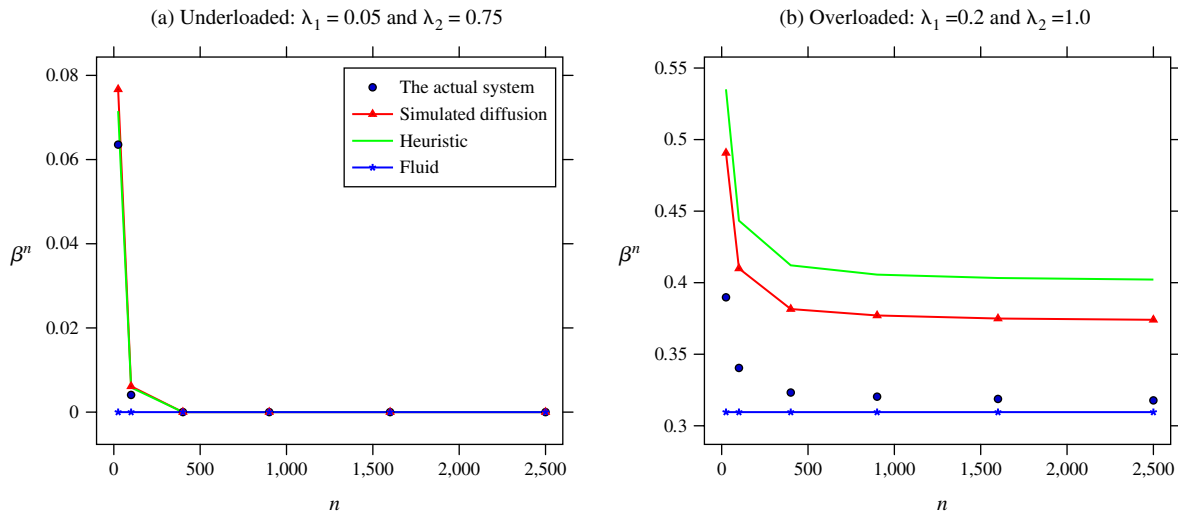
where $\Phi(\cdot)$ is the cumulative distribution function of the standard normal distribution. By (17)–(19), we can obtain C_2 , C_q , and

$$C_\beta = \frac{\Phi'(-\bar{\lambda}_1/\mu_1 - \bar{\lambda}_2/(p\mu) - (1-\phi)C_\beta)}{\Phi(-\bar{\lambda}_1/\mu_1 - \bar{\lambda}_2/(p\mu) - (1-\phi)C_\beta)}. \quad (20)$$

By (16), the delay probability of the n th system can be approximated by C_β/\sqrt{n} ; the throughput rate can be approximated by $p\mu \mathbb{E}[Z_2^n(\infty)/n] = p\mu(1 + C_2/\sqrt{n})$. Thus, the accuracy of the estimation of the system performance is reflected by the coefficients C_β and C_2 .

6.2.2. Accuracy of the Heuristic. To show how the above heuristic approximates the delay probability and

Figure 7. (Color online) Underload and Overload: The Delay Probabilities as a Function of n When $\lambda_1^n = \lambda_1 n$ and $\lambda_2^n = \lambda_2 n$



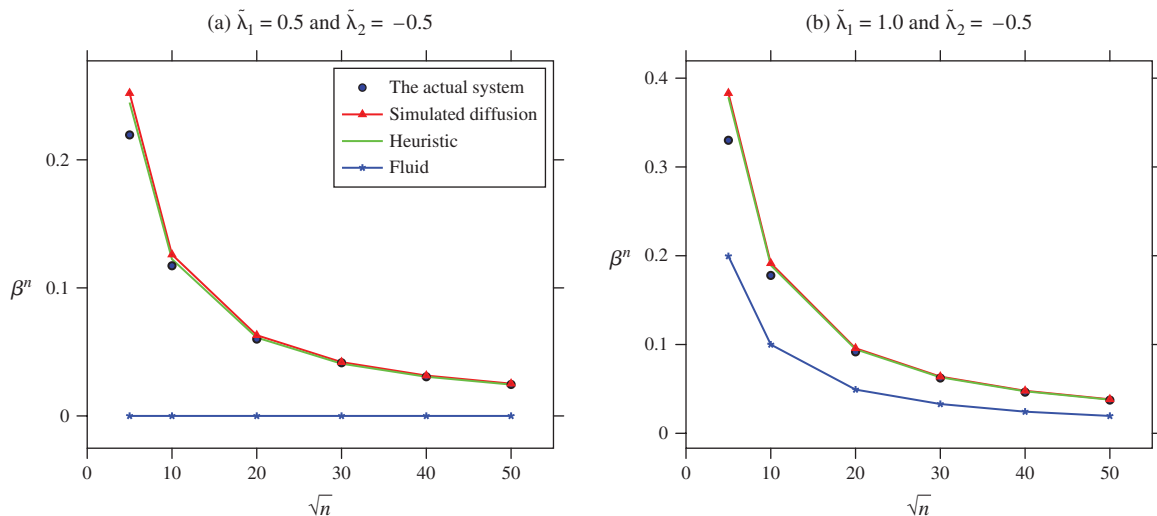
throughput rate of the diffusion-scaled processes as well as the actual system, we conduct a numerical experiment. We simulate large systems to obtain the diffusion limits and compare them with the heuristic ones. For $1/\mu_1 = 1/\mu_2 = 1$, $1/\mu_s = 0.0001$, $\theta = 0.4$, $\phi = 0.5$, and $1/\mu_f = \infty$, Figures 7 and 8 compare the diffusion-scaled delay probabilities with the heuristic ones. As expected, our heuristic mimics the performance of the simulated diffusion limits well, especially when the systems are underloaded or critically loaded. Note that, in the network, all the channels are pooled to serve both licensed and unlicensed users, while the heuristic estimates the coefficients pretending that the licensed users occupy a fixed number of channels. When the system is underloaded or critically loaded, the heuristic works well as long as the channels are well allocated to the two types of users, as shown in Figures 7(a) and 8.

The impact of decoupling the channels is higher when the system is overloaded, as shown in Figure 7(b).

6.2.3. Comparison to the Fluid Approximation. Figures 7 and 8 also plot the delay probabilities under the fluid-scaled processes and of the actual systems. As one can see, the fluid approximation always underestimates β , while the diffusion approximation always overestimates it. Furthermore, the fluid approximation outperforms the diffusion approximation when the system is overloaded, and the converse is true when the system is underloaded or critically loaded. Thus, neither method is uniformly more accurate than the other. However, further comparisons reveal the following:

1. The fluid scaling leads to analytical closed-form approximations, while analysis under the diffusion

Figure 8. (Color online) Critical Load: The Delay Probabilities as a Function of \sqrt{n} When $\lambda_1^n = \tilde{\lambda}_1 \sqrt{n}$ and $\lambda_2^n = \rho \mu n + \tilde{\lambda}_2 \sqrt{n}$



scaling involves solving the steady states of multidimensional diffusion processes, which is known to be an open question in most cases.

2. The closed-form approximations under the fluid scaling reveal important operational insights (in Section 4) that are not obvious under the diffusion approximation.

3. The diffusion approximation may not be feasible when a system is overloaded or the number of licensed users is comparable to or more than that of unlicensed ones, while the fluid approximation can be used under any load level with any ratio between the licensed and unlicensed users. Such a drawback may further limit the diffusion approximation to be adopted to systems with time-varying or random arrivals.

7. Conclusions and Future Research

Opportunistic access of licensed spectrum by unlicensed users is widely considered a way to alleviate artificial scarcity of radio spectrum by increasing the spectrum utilization. However, it may reduce the service quality for licensed users because of potential interference from unlicensed users. While much research on spectrum sharing has been conducted by researchers in electrical engineering, with the main focus on technological issues, the operational aspects have not been adequately addressed through analytical work.

In this paper, we model a shared network consisting of both licensed users and unlicensed users as a multiclass, many-server queueing system. The distinctive features of our model are that the service requirement of an unlicensed user can be fulfilled even after multiple interruptions and the unlicensed users waiting in the queue are required to sense channel availability periodically while waiting. These features complicate system dynamics and lead to quite different insights from those derived from most service systems. We show that the sensing frequency of the unlicensed users waiting in the queue does not affect system performance from the operational perspective, and its decision should be based on technological concerns. When the system is underloaded or critically loaded, there is no need to restrict the service session of unlicensed users. Otherwise, limiting the transmission of the unlicensed users is necessary only when the system load is above a threshold. Thus, it is possible to improve spectrum utilization while guaranteeing a very high service level, as expected by licensed users in practice, and spectrum sharing can potentially be a socially optimal solution to alleviating spectrum scarcity.

Spectrum sharing, if feasible, is especially beneficial for systems with a smaller portion or a large number of licensed users with shorter service times. Our study sheds light on the implementation of spectrum sharing and opens the door for new applications of

existing queueing theory in wireless communication networks, which may lead to the development of new methodologies.

Our study also provides some rich research opportunities. For instance, the arrival rates of the users may be uncertain in practice. Our preliminary result shows that higher variance will always hurt system performance if the system is expected to be underloaded or critically loaded. However, if the system is expected to be overloaded, it seems that increasing the variability up to a certain level will actually improve the throughput rate. Thus, research needs to be done to investigate the impact of uncertain arrival rates on system performance.

In reality, users' behavior in data transmission can be more complicated than those in the network in Figure 1. For instance, unlicensed users who have to abandon the system earlier may reenter the system later while licensed users may abandon the system if no idle channel is available upon arrival. Also, sensing may not be perfect (e.g., a false alarm can occur), in which case a spectrum opportunity is overlooked by an unlicensed user. It will be interesting to incorporate these elements into the model and examine how they change the system performance and operational decisions.

The insights revealed in this research may also pave the way for studying other important business issues in wireless communications, such as contract design and pricing in shared networks. For instance, how should a spectrum owner set the prices and decide the service quality to both licensed and unlicensed users in a shared network? Should unlicensed users be charged a fixed and/or usage-based fee? Since unlicensed users may belong to different service providers, should a spectrum owner run an auction to select the service providers and settle the prices?

Acknowledgments

The authors acknowledge the advice of Professor Qian Zhang from the Department of Computer Science at Hong Kong University of Science and Technology and the encouragement of the area editor, the associate editor, and the two anonymous reviewers.

References

- Akyildiz IF, Lee W-Y, Vuran MC, Mohanty S (2006) Next generation/dynamic spectrum access/cognitive radio wireless networks: A survey. *Comput. Networks* 50(13):2127–2159.
- Artalejo J, Gómez-Corral A (2008) *Retrial Queueing Systems: A Computational Approach* (Springer, Berlin).
- Atar R, Mandelbaum A, Reiman MI (2004) Scheduling a multi class queue with many exponential servers: Asymptotic optimality in heavy traffic. *Ann. Appl. Probab.* 14(3):1084–1134.
- Bansal G, Hossain J, Bhargava V (2008) Optimal and suboptimal power allocation schemes for OFDM-based cognitive radio systems. *IEEE Trans. Wireless Comm.* 7(11):4710–4718.
- Biglieri E, Goldsmith A, Greenstein L, Mandayam N, Poor H (2013) *Principles of Cognitive Radio* (Cambridge University Press, New York).

- Brosh I (1969) Preemptive priority assignment in multichannel systems. *Oper. Res.* 17(3):526–535.
- Buzen JP, Bondi AB (1983) The response times of priority classes under preemptive resume in $M/M/m$ queues. *Oper. Res.* 31(3):456–465.
- Capar F, Martoyo I, Weiss T, Jondral F (2002) Comparison of bandwidth utilization for controlled and uncontrolled channel assignment in a spectrum pooling system. *IEEE 55th Vehicular Tech. Conf.*, Vol. 3 (IEEE, Piscataway, NJ), 1069–1073.
- de Véricourt F, Jennings OB (2008) Dimensioning large-scale membership services. *Oper. Res.* 56(1):173–187.
- de Véricourt F, Zhou Y-P (2005) Managing response time in a call-routing problem with service failure. *Oper. Res.* 53(6):968–981.
- Devroye N, Mitran P, Tarokh V (2006) Achievable rates in cognitive radio channels. *IEEE Trans. Inform. Theory* 52(5):1813–1827.
- El Gamal A, Mammen J, Prabhakar B, Shah D (2006) Optimal throughput-delay scaling in wireless networks—Part I: The fluid model. *IEEE Trans. Inform. Theory* 52(6):2568–2592.
- Falin GI, Templeton JGC (1997) *Retrial Queues*, Vol. 75 (CRC Press, Boca Raton, FL).
- Goldsmith A, Jafar S, Maric I, Srinivasa S (2009) Breaking spectrum gridlock with cognitive radios: An information theoretic perspective. *Proc. IEEE* 97(5):894–914.
- Gong S, Wang P, Duan L (2015) Distributed power control with robust protection for PUs in cognitive radio networks. *IEEE Trans. Wireless Comm.* 14(6):3247–3258.
- Gupta P, Kumar P (2000) The capacity of wireless networks. *IEEE Trans. Inform. Theory* 46(2):388–404.
- Gurvich I, Perry O (2012) Overflow networks: Approximations and implications to call center outsourcing. *Oper. Res.* 60(4):996–1009.
- Gurvich I, Whitt W (2009) Queue-and-idleness-ratio controls in many-server service systems. *Math. Oper. Res.* 34(2):363–396.
- Halfin S, Whitt W (1981) Heavy-traffic limits for queues with many exponential servers. *Oper. Res.* 29(3):567–588.
- Hossain E, Niyato D, Han Z (2009) *Dynamic Spectrum Access and Management in Cognitive Radio Networks* (Cambridge University Press, Cambridge, UK).
- Huang S, Liu X, Ding Z (2008) Opportunistic spectrum access in cognitive radio networks. *Proc. 27th Conf. Comp. Comm.* (IEEE INFOCOM 2008) (IEEE, Piscataway, NJ), 2101–2109.
- Hunt PJ, Kurtz TG (1994) Large loss networks. *Stochastic Processes Their Appl.* 53(2):363–378.
- Jaiswal NK (1968) *Priority Queues*, Mathematics in Science and Engineering (Elsevier Science, New York).
- Jia J, Zhang Q, Shen X (2008) HC-MAC: A hardware-constrained cognitive mac for efficient spectrum management. *IEEE J. Selected Areas Comm.* 26(1):106–117.
- Kurtz TG (1992) Averaging for martingale problems and stochastic approximation. *Applied Stochastic Analysis*, Lecture Notes in Control and Information Sciences, Vol. 177 (Springer, Berlin), 186–209.
- Liu KR, Wang B (2010) *Cognitive Radio Networking and Security: A Game-Theoretic View* (Cambridge University Press, New York).
- Luo J, Zhang J (2013) Staffing and control of instant messaging contact centers. *Oper. Res.* 61(2):328–343.
- Maglaras C, Zeevi A (2004) Diffusion approximations for a multiclass Markovian service system with “guaranteed” and “best-effort” service levels. *Math. Oper. Res.* 29(4):786–813.
- Maglaras C, Zeevi A (2005) Pricing and design of differentiated services: Approximate analysis and structural insights. *Oper. Res.* 53(2):242–262.
- Mandelbaum A, Massey W, Reiman M, Stolyar A, Rider B (2002) Queue lengths and waiting times for multiserver queues with abandonment and retrials. *Telecomm. Systems* 21(2–4):149–171.
- Mishra S, Sahai A, Brodersen R (2006) Cooperative sensing among cognitive radios. *Proc. IEEE Internat. Conf. Comm.*, Vol. 4 (IEEE, Piscataway, NJ), 1658–1663.
- Mitola J, Maguire JGQ (1999) Cognitive radio: Making software radios more personal. *IEEE Personal Comm.* 6(4):13–18.
- Pang G, Perry O (2015) A logarithmic safety staffing rule for contact centers with call blending. *Management Sci.* 61(1):73–91.
- Perry O, Whitt W (2011a) A fluid approximation for service systems responding to unexpected overloads. *Oper. Res.* 59(5):1159–1170.
- Perry O, Whitt W (2011b) An ODE for an overloaded X model involving a stochastic averaging principle. *Stochastic Systems* 1(1):59–108.
- Perry O, Whitt W (2013) A fluid limit for an overloaded X model via an averaging principle. *Math. Oper. Res.* 38(2):294–349.
- Prasad R, Dixit S, van Nee R, Ojanpera T, eds. (2010) *Globalization of Mobile and Wireless Communications: Today and in 2020*, Signals and Communication Technology (Springer, Dordrecht, Netherlands).
- Puhalskii AA, Reiman MI (2000) The multiclass $GI/PH/N$ queue in the Halfin-Whitt regime. *Adv. Appl. Probab.* 32(2):564–595.
- van de Beek J, Riihijarvi J, Achtzehn A, Mahonen P (2012) TV white space in Europe. *IEEE Trans. Mobile Comput.* 11(2):178–188.
- Wang J, Baron O, Scheller-Wolf A (2015) $M/M/c$ queue with two priority classes. *Oper. Res.* 63(3):733–749.
- White H, Christie LS (1958) Queuing with preemptive priorities or with breakdown. *Oper. Res.* 6(1):79–95.
- Whitt W (2002) *Stochastic-Process Limits*. Springer Series in Operations Research (Springer-Verlag, New York).
- Whitt W (2006) Fluid models for multiserver queues with abandonments. *Oper. Res.* 54(1):37–54.
- Yang T, Templeton J (1987) A survey on retrial queues. *Queueing Syst.* 2(3):201–233.
- Yom-Tov G, Mandelbaum A (2014) Erlang-R: A time-varying queue with reentrant customers, in support of healthcare staffing. *Manufacturing Service Oper. Management* 16(2):283–299.
- Zhan D, Ward AR (2014) Threshold routing to trade off waiting and call resolution in call centers. *Manufacturing Service Oper. Management* 16(2):220–237.
- Zhang Q, Jia J, Zhang J (2009) Cooperative relay to improve diversity in cognitive radio networks. *IEEE Comm. Magazine* 47(2):111–117.
- Zhao Q, Geirhofer S, Tong L, Sadler B (2008) Opportunistic spectrum access via periodic channel sensing. *IEEE Trans. Signal Processing* 56(2):785–796.

Shining Wu is an assistant professor in logistics and maritime studies at the Hong Kong Polytechnic University. His research interests include supply chain management, strategic consumer behavior, queueing theory and its applications, operational issues in spectrum sharing, and data-driven optimization for queueing systems.

Jiheng Zhang is an associate professor in industrial engineering and decision analytics at the Hong Kong University of Science and Technology. His research interests are in applied probability, stochastic modeling and optimization, data analysis, numerical methods, and algorithms.

Rachel Q. Zhang is a professor in industrial engineering and decision analytics at the Hong Kong University of Science and Technology. Her research interests include supply chain and inventory management, stochastic analysis of service operations, and the interface of finance and operations.