

## **Title: Spectrum Access Games and Strategic Learning in Cognitive Radio Networks for Delay-Critical Applications**

**Submission by/Based on: UCLA tutorial** - van der Schaar & Fu based on paper, "Learning to Compete for Resources in Wireless Stochastic Games " by F. Fu and M. van der Schaar (Trans. On Vehicular Technology)

**Special issue submission date:** 28 December 2007

**Abstract:** With the current proliferation of high bandwidth and delay-sensitive multimedia applications and services, each wireless user will try to maximize its utility by acquiring as much spectrum resources as possible, unless a preemptive mechanism exists in the network. Thus, emerging solutions for dynamic spectrum access in cognitive radio networks will need to adopt market-based approaches in order to effectively regulate the available resources. In this paper, we show how various centralized and decentralized spectrum access markets can be designed based on a stochastic game framework, where multimedia users (also referred to as secondary users) can compete over time for the dynamically available transmission opportunities (spectrum "holes"). When operating in such spectrum access markets, wireless users become selfish, autonomous agents that strategically interact in order to acquire the necessary spectrum opportunities. We also show how wireless users can successfully compete with each other for the limited and time-varying spectrum opportunities, given the experienced dynamics in the wireless network, by optimizing both their external actions (e.g. the resource bids, power and channel used for transmission etc.) and internal actions (e.g. the modulation schemes etc.). To determine their optimal actions in an informationally-decentralized setting, users will need to learn and model directly or indirectly the other users' responses to their external actions. We study the outcome of various dynamic interactions among self-interested wireless users possessing different knowledge, and determine that the proposed framework can lead to multi-user communication systems that achieve new measures of optimality, rationality and fairness. Finally, the illustrative results show that the presented game-theoretic solution for wireless resource management enables users that deploy enhanced (smarter) learning and communication algorithms to derive higher utilities.

Formatted: Highlight

### **1. Introduction**

#### **1.1. Motivation**

Due to their flexible and low cost infrastructure, wireless networks are poised to enable a variety of delay-sensitive multimedia transmission applications, such as videoconferencing, emergency services, surveillance, telemedicine, remote teaching and training, augmented reality, and distributed gaming. However, existing wireless networks provide limited, time-varying resources with only limited support for the Quality of Service (QoS) required by the *delay-sensitive, bandwidth-intense and loss-tolerant multimedia applications*. The scarcity and variability of resources does not significantly impact delay-insensitive applications (e.g., file transfers), but has considerable consequences for multimedia applications and often leads to unsatisfactory user experience.

In recent years, the research focus has been to adapt the resource allocation (e.g. integrated and differentiated services) methods, and transmission (e.g. TCP) strategies and concepts designed for the wired (Internet, ATM) communications to the time-varying and bandwidth-constrained wireless networks. However, such solutions (e.g. QoS enabled 802.11e solutions) do not provide fair or efficient support for delay-sensitive applications such as multimedia streaming in crowded or dynamic wireless networks [13], because they ignore the wireless system dynamics, including the time-varying source and channel characteristics, the mobility of the wireless sources, the unpredictability of wireless users or interference sources coming or leaving the network, etc.

One vision for emerging cognitive radio networks assumes that certain portions of the spectrum will be opened up for secondary users<sup>1</sup> (SUs), such as wireless multimedia applications, to autonomously and opportunistically share the spectrum becoming available once primary users (PUs) are not active [1][2][3]. Importantly, in cognitive radio networks, heterogeneous wireless users (with different utility-rate functions, delay tolerances, traffic characteristics, knowledge and adaptation abilities) will need to coexist and interact within the same band [6]. However, to enable the proliferation of multimedia applications over cognitive radio networks, solutions for dynamic spectrum access will need to consider different challenges that are discussed next.

## ***1.2. Challenges for Dynamic Spectrum Access in Cognitive Radio Networks***

Next generation networking solutions will need to address, besides other issues related to the co-existence between the PUs and SUs, the following four challenges associated with designing efficient resource

---

<sup>1</sup> The secondary users/applications are envisioned in this paper to be a single transmitter-receiver pair.

management solutions for dynamic and autonomous applications over wireless environments.

- A first challenge arises due to the **dynamic, time-varying nature of applications, source and channel characteristics**. As the source characteristics are changing, the delays that are tolerable at the application layer and the derived utility (e.g. quality or fidelity) can vary significantly. This influences the performance of the different transmission strategies at the various layers and, ultimately, the choice of the optimal strategy adopted by the transmitter. Hence, the utility that a user derives from using a certain resource dynamically varies over time, depending on both the “environment” (e.g. application, source and channel characteristics), which is not in the control of the user, as well as on the user’s response to this environment, which is the selected transmission strategy (at the application, transport, network, MAC or physical layers).
- A second challenge associated with multi-user transmission and resource management is that the **wireless users’ actions and their performances are coupled** [2][43], since the transmission strategy of a user impacts and is impacted by the competing users (see Figure 1). Hence, a user’s actions will have a direct impact not only on its own utility, but also on the performance of the other wireless users sharing the spectrum.
- A third challenge comes from the **informationally-decentralized nature of the multi-user wireless resource management problem**. Each wireless user can derive different utilities based on the amount of resources consumed/allocated, which also depends on its “private” information (i.e. traffic and channel characteristics, and selected transmission strategy). However, in general, in a practical transmission scenario, the private information of each user is not known by the resource manager or other wireless users (see Figure 1). Moreover, the users are not always directly aware of the resources requested by other users, how the other wireless users allocate their power etc. This is different from the multi-access channel setting in conventional communication systems, where the assumption is made that the private information of all users is known by the resource coordinator in the centralized network setting. Note also that since this information is private, wireless users may lie when declaring their private information.
- Finally, most existing wireless resource management solutions disregard two important properties of the **autonomous wireless users**: their knowledge (and thus ability to learn and optimize their transmission strategies by anticipating the coupling with the other users and the impact of their



architecting communications systems that is governed by dynamic spectrum access rules [10], and where SUs can compete with each other based on their available transmission strategies as well as their knowledge about the environment and other SUs. Specifically, wireless users can compete for spectrum access based on “market” rules designed using either non-collaborative (e.g. mechanism design [15][22]) or cooperative<sup>2</sup> (e.g. bargaining or coalition theory [33]) game-theory (see Figure 2). Hence, wireless users become rational players competing for the available wireless resources in a stochastic (or repeated) game, played repeatedly by the communication system entities. To maximize their utilities, the users will need to negotiate or acquire spectrum access as well as to proactively adapt their cross-layer transmission strategies. Note that the competition is performed using incomplete information about other users’ private information, actions or utility functions, and it is influenced by the SUs’ behaviors – e.g. attitudes towards risk, willingness to pay for resources, maliciousness etc.

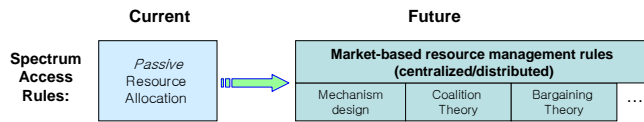


Figure 2. Evolution of spectrum access rules to create a dynamic wireless resource market.

The paradigm proposed in this paper is referred to as *knowledge-driven networking*, since the various network entities (spectrum moderators, access points or wireless users) will need to make their decisions about spectrum division, spectrum negotiation, or cross-layer optimization based on their knowledge about the environment and other network entities (PUs and SUs). As mentioned previously, since the decisions that need to be taken are based on *incomplete information* about the environment and other entities (see Figure 3), the knowledge that a network entity possesses will influence its efficiency and performance. By gathering information (private observations or explicit information exchanges with other SUs), and subsequently learning and reasoning based on this information, network entities can develop true beliefs<sup>3</sup> about the current state of the communication system and its evolution over time, based on which they can select the optimal policies for interacting with other entities such that they maximize their utilities. For example, a foresighted user can learn the other entities’ responses to its actions, thereby being able to forecast the impact of its actions on the wireless system and, ultimately, to optimize its

<sup>2</sup> Cooperative game theory is a parallel branch to the more widely known topic of non-cooperative game theory. The term “cooperative” does *not* mean that users have interests that are aligned, but rather cooperative game theory concepts are relevant in situations where a scarce resource is to be divided *fairly* among competing users. Concepts, such as the bargaining solutions, embody specific notions of fairness and take into account the strategic interests of competing users [33].

<sup>3</sup> Plato defined knowledge as “justified true belief” [<http://en.wikipedia.org/wiki/knowledge>].

resulting utility over time, rather than just myopically optimize its immediate performance.

Summarizing, in this paper, we show how users can compete for resources in various wireless markets and briefly introduce the necessary principles and methods for:

- designing different dynamic spectrum access rules for a variety of communication scenarios;
- enabling the wireless users to learn the system dynamics based on observations and/or explicit information exchanges, and improve their strategies for playing the spectrum access game;
- evaluating the “value” of learning, and “value of information” for a user in terms of its utility impact;
- coupling the internal and external actions<sup>4</sup> of the wireless users to allow them to achieve an optimal response to the dynamically changing wireless resource market.

Our main focus in this paper will be on designing solutions for emerging cognitive radio networks in which wireless stations<sup>5</sup> (WSTAs) are able to utilize multiple frequency bands, thereby allowing WSTAs to dynamically harvest additional resources. However, the proposed solutions will also be beneficial when deployed in existing ISM radio bands, dedicated bands, or in first-generation versions of cognitive radio networks, which may only rely in their implementations on multiple ISM bands. Thus, the focus of this paper will be on designing new dynamic spectrum access and strategic transmission solutions, but not on detecting primary users and identifying spectrum opportunities for WSTAs. For this topic, we refer the interested reader to [4][5] as well as to several articles in this special issue, which are addressing these important issues. In this paper, we assume that the spectrum opportunities can be known by simply accessing a dynamically created Spectrum Opportunity Map [11].

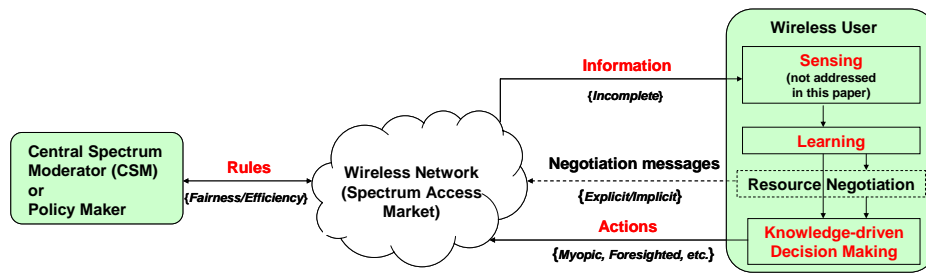


Figure 3. Knowledge-driven wireless networking.

#### 1.4. Paper layout

<sup>4</sup> The internal and external actions and the coupling between them will be discussed in detail in Section 3.

<sup>5</sup> In this paper, the secondary users (SUs) are also termed wireless stations (WSTAs) and thus, we use the denominations “SU” and “WSTA” interchangeably.

The paper is organized as follows. First, we discuss in Section 2 several important issues that need to be considered when defining the spectrum access game among the wireless users. In Section 3 we define a general framework for constructing the wireless resource game. Section 4 presents learning solutions that can be deployed by the users to improve their performance when playing the spectrum access game. This section also presents several illustrative examples for the presented framework. However, note that these results are only illustrative and a significant body of work will need to take place before comprehensive solutions can be implemented based on the presented knowledge-driven networking framework. Section 5 concludes the paper by highlighting the impact of developing such a knowledge-driven networking framework.

## 2. Issues to consider for the design and construction of spectrum access games

In this section, we summarize several key issues that need to be considered when designing and constructing any spectrum access game (“market”) for cognitive radio networks.

- **Resource types**

In [46], the resources in a certain market were classified according to different criteria – continuous vs. discrete, divisible or not, sharable or not, static or not. These classification criteria are also useful and should be considered by both the WSTAs and spectrum regulators when accessing or determining the spectrum division rules and policies for cognitive radio networks. For instance, the resources can be discrete (e.g. in FDMA, where a single channel is allocated to each user, or even in the case of TDMA, when users are allocated a certain percentage of a service time interval to access the channel), or they can be continuous (e.g. the adjustment of power levels). The resources can be defined as sharable or not, depending on the wireless protocol used. For instance, in FDMA, only one user can share a frequency band, and in TDMA, multiple users can time-share the same channel access. The resources are static or dynamically varying over time, e.g. depending on the PU access.

- **Stochastic vs. repeated games**

Stochastic games [18] are dynamic, competitive games with probabilistic transitions played by several SUs. The game is played in a sequence of stages. At the beginning of each stage, the game is in a certain state. The SUs select their actions and each SU receives a reward that depends on both its current state

Deleted: wireless stations<sup>6</sup> (

Deleted: )

and its selected external and internal actions. The game then moves to a new state with a certain probability, which depends on the previous state and the actions chosen by the SUs. This procedure is repeated at the new state and the interaction continues for a finite or infinite number of stages. The stochastic games are generalizations of repeated games, which correspond to the special case where there is only one state.

- **One shot vs. multi-stage games**

The games taking place in the wireless networks can be categorized as one-shot or multi-stage games, depending on whether the allocation is performed once or repeatedly. For instance, in 802.11e, the resource allocation is usually performed by the wireless access point only once, when a SU joins the network [13]. The advantages of such one-shot allocations are that the complexity associated with implementing any resource allocation is kept limited. However, the disadvantage is that this solution does not consider the time-varying source and channel characteristics of the SUs, and the static allocation may become inefficient over time [13][22]. In this case, repeated or stochastic games can be defined, where the users repeatedly compete for the available resources at each stage of the wireless resource allocation game.

- **Centralized vs. decentralized**

In the centralized setting, a central spectrum moderator (CSM) such as an access point or base station is responsible for determining and enforcing the allocation among the competing users. In the decentralized setting, the SUs interact with each other directly, through the actions they perform [34][35], and there is no moderator involved in the negotiation. Note that in current ISM bands, the wireless users are using the same spectrum access protocols and thus, distributed solutions can be easily designed and enforced. However, in cognitive radio networks, the SUs will be heterogeneous in terms of protocols, utility-cost functions etc., and this needs to be explicitly considered when designing distributed solutions/protocols.

- **Budget-balanced vs. money-making resource allocation solutions**

Wireless resource allocation solutions can be budget balanced (i.e. all the money users pay to the wireless network is allocated back to them) or money-making. For instance, many well-known mechanism implementations, such as the Vickrey-Clarke-Groves (VCG) mechanism [25], charge users for using resources. In the VCG mechanism, the transfer (money, tokens etc.) is only delivered from the direction of the SU to the moderator, which becomes a money-making (transfer-making) entity. The moderator can



use this transfer to maintain or upgrade its service, purchase additional spectrum etc. However, if there is either no moderator in the system or the participating SUs want to prevent the moderator from behaving as a profit maker (e.g. in some wireless LAN usage scenarios), which may potentially result in the moderator trying to alter the users' allocations in order to maximize its revenue, the SU will need to deploy a budget-balanced mechanism (i.e. all the money users pay to the wireless network is allocated back to them) [15]. Note also that the wireless resource markets can also be designed and regulated without transfers. However, such transfers provide important benefits for wireless networks, because they force multimedia users to truthfully reveal their resource needs [22].

- **Social decisions (fairness rules)**

Various fairness rules can be imposed by the CSM or can be negotiated in a decentralized manner by the WSTAs, e.g. using bargaining solutions. Some examples investigated in current wireless networks are weighted-sum maximizations of rates or utilities among the participating users, envy-free fairness solutions or egalitarian solutions. For a comprehensive discussion of these fairness rules, the interested reader is referred to [46]. Performing the resource allocation in the utility domain rather than the resource domain is vital for multimedia users and can result in significant performance gain over application-agnostic resource allocation solutions [31].

- **Desired equilibrium concepts**

When playing or designing wireless resource games, SUs or moderators will need to *proactively* negotiate, select or design their desired equilibrium point. This is unlike most game-theory literature [23], which is developing descriptive models (in e.g. social or biological systems) to show that certain equilibrium exist. In wireless communication games, *constructive* models are required, where equilibrium can be designed or influenced by the participating SUs (e.g. [30]). Depending on whether the game is centralized or decentralized, the CSM or the WSTAs may strive towards implementing such equilibriums. Note that often other equilibrium concepts rather than the well-known Nash equilibrium are desired. Examples are correlated equilibriums [40][41], dominant strategy equilibriums [22], Stackelberg equilibriums etc. For instance, in [30], to characterize the multi-user interaction in the distributed power-control game where a foresighted SU can anticipate the responses of its opponent SUs to its actions, the Stackelberg equilibrium is introduced which is shown to outperform the well-known Nash equilibrium.

- **Implementation complexity**

An important issue associated with the implementation and adoption of wireless resource markets is the resulting complexity for both the CSM, which needs to implement the different resource allocations, and the SUs, which may adopt strategic learning algorithms to be able to compete against other SUs. Hence, new metrics such as the value of learning or the value of information exchanges (which will be discussed in Section 4) need to be deployed to trade-off the actual benefit that the network entities can derive by increasing their knowledge against the expense of a higher complexity cost.

### 3. Dynamic multi-user spectrum access games

While the knowledge-driven framework presented in this paper can be implemented in most network settings discussed in the previous section, we will illustrate in this paper only several specific wireless transmission scenarios. In particular, as mentioned in the introduction, we focus on developing wireless resource markets for secondary networks (SN). In SN, the secondary users (SUs) can opportunistically utilize the network resources that are vacated by the PUs. For illustration purposes, we assume that the SN consists of  $M$  SUs, which are indexed by  $i \in \{1, \dots, M\}$ . The SUs compete for the dynamically available transmission opportunities based on their own “private” information, knowledge about other WSTAs and available resources (and/or PUs’ behaviors). In each time slot  $\Delta T$ , the WSTAs compete with each other for spectrum access and, given the allocated transmission opportunities<sup>7</sup>, they deploy optimized cross-layer strategies to transmit their delay-sensitive bitstreams.

During each time slot, a “state” of network resources can be defined to represent the available transmission opportunities in a SN, which is denoted by  $w \in \mathcal{W}$ , where  $\mathcal{W}$  is the set of possible resource states. We can also interpret the state of network resources to reflect the behaviors of PUs in the cognitive radio networks [11]. We can also define “states” for the WSTAs. For instance, the states may represent their private information, which includes the traffic and channel characteristics. The current state of a WSTA  $i$  is denoted by  $s_i \in \mathcal{S}_i$ , where  $\mathcal{S}_i$  is the set of possible states of WSTA  $i$ .

At each time slot, WSTA  $i$  will deploy an action to compete for the network resources. This action is referred to as the external action denoted by  $a_i \in \mathcal{A}_i$ , where  $\mathcal{A}_i$  is the set of possible external actions. An example of external actions in wireless networks is the selected transmit power in interference channels or the declared resource request like the TSPEC in 802.11e WLANs. Besides the external action, WSTA  $i$  will also deploy an internal action in order to transmit the delay-sensitive data. The internal action can be

---

<sup>7</sup> Note that the resource competition and data transmission may take place concurrently.

an action profile including all or a subset of actions from different layers (e.g. adaptation of the packet scheduling strategy, which error correcting codes or retransmission limits to use etc.). This action is denoted as the internal action denoted by  $b_i \in \mathcal{B}_i$ , where  $\mathcal{B}_i$  is the set of possible internal actions. Note that the external and internal action selections are coupled together as shown in [22]. Moreover, the actions' adaptation can be driven by cross-layer optimization.

In this paper, we formulate the multi-user wireless resource competition as a stochastic game<sup>8</sup>. Formally, the stochastic game is defined as a tuple  $(\mathcal{I}, \mathcal{S}, \mathcal{W}, \mathcal{A}, \mathcal{B}, P_s, P_w, \mathcal{R})$ , where  $\mathcal{I}$  is the set of agents (SUs), i.e.  $\mathcal{I} = \{1, \dots, M\}$ ,  $\mathcal{S}$  is the set of state profiles of all SUs, i.e.  $\mathcal{S} = \mathcal{S}_1 \times \dots \times \mathcal{S}_M$  with  $\mathcal{S}_i$  being the state set of SU  $i$ ,  $\mathcal{W}$  is the set of network resource state.  $\mathcal{A}$  is the joint external action space  $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_M$ , with  $\mathcal{A}_i$ <sup>9</sup> being the external action set available for SU  $i$  to play the resource sharing game, and  $\mathcal{B}$  is the joint internal action space  $\mathcal{B} = \mathcal{B}_1 \times \dots \times \mathcal{B}_M$ , with  $\mathcal{B}_i$  being the internal action set available for SU  $i$  to transmit delay-sensitive data.  $P_s$  is a transition probability function defined as a mapping from the current state profile  $s \in \mathcal{S}$ , corresponding joint external actions  $a \in \mathcal{A}$  and internal actions  $b \in \mathcal{B}$  and the next state profile  $s' \in \mathcal{S}$  to a real number between 0 and 1, i.e.  $P: \mathcal{S} \times \mathcal{A} \times \mathcal{B} \times \mathcal{S} \mapsto [0, 1]$ .  $P_w$  is a transition probability function defined as a mapping from the current resource state  $w \in \mathcal{W}$  and the next state  $w' \in \mathcal{W}$  to a real number between 0 and 1, i.e.  $P: \mathcal{W} \times \mathcal{W} \mapsto [0, 1]$ . This will be discussed subsequently, in more detail.  $\mathcal{R}$  is a reward vector function defined as a mapping from the current state profile  $s \in \mathcal{S}$  and corresponding joint external and internal actions  $a \in \mathcal{A}$  and  $b \in \mathcal{B}$  to an  $M$ -dimensional real vector with each element being the reward to a particular agent, i.e.  $\mathcal{R}: \mathcal{S} \times \mathcal{A} \times \mathcal{B} \mapsto \mathbb{R}^M$ .

In the cognitive radio environment, if the secondary network shares all the spectrum resources with the primary network, it could happen that all the resources are occupied by the primary users and thus, no "spectrum holes" are available for the secondary users. In this case, no QoS guarantee can be provided to the secondary users. However, in most usage scenarios, not all primary users are active simultaneously and even if they are all active simultaneously, they will not use the spectrum continuously. Thus, most of the time, even when all the primary users are active, the secondary network will have access to limited resources, which can be used to guarantee the minimum resource needs of the secondary users. Thus, in

<sup>8</sup> The use of games for dynamic spectrum access in cognitive radio networks were discussed already in e.g. [2][32].

<sup>9</sup> Note that the action set may depend on the state of the SU. For simplicity, we assume that the actions sets are the same for all the states of the SU.

Deleted: generalized

Deleted: generalized

Formatted: Font: Not Italic

this paper, we assume that the secondary network has always access to a limited amount of resources. It should be noted that this assumption does not violate the stochastic game model for the secondary network.

The state transition for the network resource state is determined by the PUs, and not by the SUs. In other words, the SUs' actions will not affect the network resource state transition. This structure actually creates an opportunity to allow the PUs to be agents with higher priorities in this stochastic game. Moreover, multiple parallel games can be easily defined in this way for different priority users of the same wireless infrastructure [50]. According to how the WSTAs compete for spectrum access and exchange information about (and access) the available spectrum opportunities, we consider two types of stochastic games for wireless resource "markets": centralized stochastic games and distributed stochastic games.

**Deleted:** The difference between this generalized wireless stochastic game proposed here and the conventional stochastic game [18] is the common state  $w \in \mathcal{W}$  that represents the common information (e.g. the spectrum opportunity map discussed previously).

**Deleted:** generalized

**Deleted:** It is easy to see that the conventional stochastic game is a special case of our definition by assuming the network resource state  $w$  is constant.

**Deleted:** generalized

**Deleted:** generalized

**Deleted:** generalized

### 3.1 Centralized stochastic game

In the centralized stochastic game, the competition between WSTAs is coordinated by a CSM, which can be an access point, base station or selected leader. Specifically, at each stage, the WSTAs perform the external actions  $a_i$  (e.g. resource requirement, competition bids) and send the CSM a message  $m_i$  representing the selected actions. An example of a wireless infrastructure where such a centralized stochastic game can be implemented are wireless LANs (802.11a PCF or 802.11e HCF), where the CSM role is played by the access point.

After receiving the messages  $\mathbf{m} = [m_1, \dots, m_M]$  from all the WSTAs, the CSM performs the resource allocation according to a certain rule, i.e.

$$[\mathbf{r}_1, \dots, \mathbf{r}_M] = f(\mathbf{m}, w), \quad (1)$$

where  $\mathbf{r}_i$  is the resource allocation to WSTA  $i$ , and  $f(\cdot, \cdot)$  represents the resource allocation rule based on the announced message  $\mathbf{m}$  and network resource state  $w$ .

After receiving the resource allocation  $\mathbf{r}_i$ , WSTA  $i$  performs its own internal action  $b_i$  to transmit the delay-sensitive data based on its current state  $s_i$ . Note that in the centralized game, the resource allocation  $\mathbf{r}_i$  for each SU  $i$  is computed by the CSM based on the external actions of all the WSTAs. The state transition can be represented by

$$s_i^+ = g_i(s_i, \mathbf{r}_i, b_i), \quad (2)$$

and the reward function is computed as

$$R_i = h_i(s_i, r_i, b_i). \quad (3)$$

The states and reward functions for an SU as well as the coupling with the other SUs will be discussed in subsequent sections.

The centralized stochastic game for the cognitive radio network is illustrated in Figure 4. This can be employed across multiple channels (frequency bands) simultaneously.

Deleted: generalized

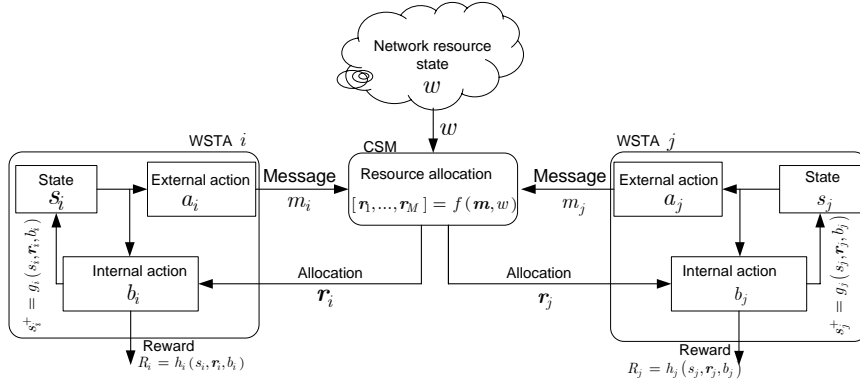


Figure 4. Message exchange between WSTAs and CSM, and the actions performed by WSTAs.

Each WSTA plays the centralized stochastic game against the other WSTAs by not only selecting its (external) actions, but also by selecting and implementing its internal actions for data transmission. Moreover, the state transition of each WSTA  $i$  is directly impacted by both its external and internal actions as well as the external actions of other WSTAs through the resource allocation. The same holds true for the reward. The cross-layer transmission strategies constitute the internal actions deployed by a WSTA. When determining its external actions, a WSTA will need to predict not only what will be the evolution of the source and channel characteristics over time, but also the cross-layer strategy that the user will select given the future environment condition. As shown in [22], the cross-layer transmission strategy will not only impact the immediate reward derived by the WSTA based on transmitting the current packets, but also the future states and rewards. This is because the current cross-layer strategy will determine which packets get transmitted and thus, what are the remaining packets to be transmitted etc., which affects the future states. Hence, as shown in [22], the ability of a WSTA to adopt more efficient transmission algorithms at the various layers as well as optimize its cross-layer transmission strategies, significantly impacts the performance of both the WSTA and that of its competing WSTAs. Moreover, the behavior of a WSTA that is risk adverse or risk taking will significantly influence the way in which both its internal and external actions are selected [22]. For instance, a risk adverse WSTA may decide to

Deleted: generalized

schedule its important packets as soon as possible, at the possible expense of a higher transmission cost or of other less important packets not being transmitted.

### 3.2 Distributed stochastic game

In the distributed stochastic game, the SUs simultaneously compete for the spectrum opportunities in the absence of a CSM that coordinates their interactions. In the distributed game, no moderator exists. Examples of such distributed games are the power control games in interference channels. (For instance, the distributed power control games in e.g. [44][49] can be represented using the stochastic game formulation presented here.) In the distributed stochastic game, the WSTAs simultaneously implement the internal and external actions. However, the interactions between WSTAs are realized through the external actions. From the perspective of each WSTA, the impact from other WSTAs is aggregated into the experienced channel interference  $e(s_{-i}, a_{-i})$ . In power control games, the external action can be the power allocation, while the internal action can be the modulation and channel coding scheme. Hence, in distributed stochastic games, the reward of each WSTA  $i$  is determined by

$$R_i = h_i(s_i, a_i, b_i, e(s_{-i}, a_{-i}), w), \quad (4)$$

where  $-i$  is the set of WSTAs except WSTA  $i$ . The state transition is determined by

$$s_i^+ = g_i(s_i, a_i, b_i, e(s_{-i}, a_{-i}), w). \quad (5)$$

The states and reward functions for an SU will be discussed in subsequent sections. The distributed stochastic game for the cognitive radio network is illustrated in Figure 5.

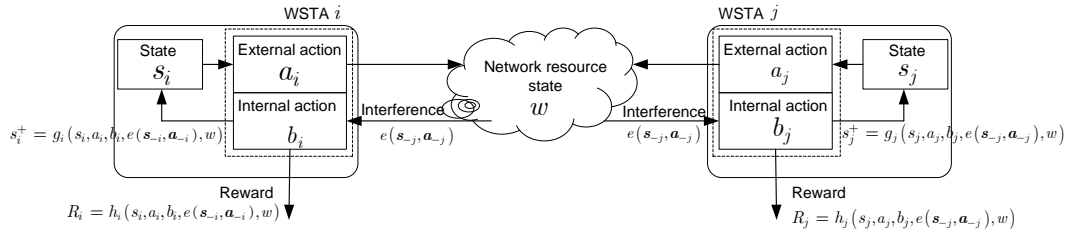


Figure 5. Actions performed by WSTAs in the distributed stochastic game.

### 3.3 Specification of the centralized stochastic game

As illustration, we consider that the SN can be formed across  $N$  channels, each indexed by  $j \in \{1, \dots, N\}$ . At each time slot, each channel is assumed to be in one of the following two states: ON (this channel is currently used by the PUs) or OFF (this channel is not used by the PUs and hence can be used by the SUs). Within each time slot, the channel is only OFF or ON [11]. At time slot  $t \in \mathbb{N}$ , the

availability of each channel  $j$  is denoted by  $w_j^t \in \{0,1\}$ , with  $w_j^t$  being 0 if the channel is in the ON state, and being 1 if it is in the OFF state. The channel availability profile for the  $N$  channels is represented by  $\mathbf{w}^t = [w_1^t, \dots, w_N^t]$ , which is the state of the network resource at time slot  $t$ . As mentioned before, this can be characterized using a Spectrum Opportunity Map [11] in the CSM.

If the CSM performs imperfect spectrum sensing as in [5], this imperfect detection only affects the common system state  $\mathbf{w}^t = [w_1^t, \dots, w_N^t]$  which is announced to all the secondary users. In this case, instead of announcing the exact common resource state, the CSM can announce the probability of the channel being available to the secondary users or not. Then, the secondary users will compete for the resource based on the probability of the channel availability. This relaxation does not provide too much insight on how the secondary users compete for the available resources with each other, but it increases the communication overhead by exchanging the probability of channel availability. Hence, in this paper, we focus on the case in which the CSM performs the perfect spectrum sensing and efficiently allocates the detected spectrum among the competing secondary users.

As in [14], we assume that a polling-based medium access protocol is deployed in the secondary network, which is arbitrated by a CSM. The polling policy is changed only at the beginning of every time slot. For simplicity, we assume that each SU can access a single channel and that each channel can be accessed by a single SU within the time slot. The SUs can switch the channels only when crossing time slots. Note that this simple medium access model used for illustration in this paper can be easily extended to more sophisticated cognitive radio models [12], where each SU can simultaneously access multiple channels or the channels are being shared by multiple SUs etc.

#### A. Wireless Stations States

We assume that WSTAs need to transmit delay sensitive applications. The bitstream at the application layer is packetized with an average packet length  $\ell$ . In this paper, we consider multimedia applications, where the application packets have a hard delay deadline, i.e. the packets will expire  $J$  stages after they are ready for transmission. Then, we can define the state of the buffer as  $\mathbf{v}_i^t = [v_{i1}^t, \dots, v_{iJ}^t]^T$ , where  $v_j^t$  ( $1 \leq j \leq J$ ) is the number of packets waiting for transmission that have a remaining life time of  $j$  time slots.

The condition of channel  $j$  experienced by WSTA  $i$  is represented by the Signal-to-Noise Ratio (SNR) and it is denoted as  $c_{ij}^t$  in dB. The channel condition profile is given by  $\mathbf{c}_i^t = [c_{i1}^t, \dots, c_{iN}^t]$ . To

model the dynamics experienced by WSTA  $i$  at time  $t$  in the cognitive radio network, we define a “state”  $s_i^t = (v_i^t, c_i^t) \in \mathcal{S}_i$ , which encapsulates the current buffer state as well as the state of each channel.

The environment experienced by each WSTA is characterized by the packet arrivals from the (multimedia) source (i.e. source/traffic characterization) connected with the transmitter, the spectrum opportunities released by Pus, and the channel conditions. Different models can be used by a WSTA to characterize the environment. However, the accuracy of the deployed models will only affect the performance of the solution, and not the general framework for multi-user interaction presented here.

### B. Internal and external actions

At the beginning of each time slot, each WSTA deploys an external action  $a_i$  to compete for the spectrum opportunities with other WSTAs. The selection of external actions will be discussed in Section 3.3.E. After receiving the resource allocation  $r_i$  from the CSM, the WSTA will deploy the internal action  $b_i^t$ . The internal action in this example includes the modulation scheme  $\gamma_i^t \in \Upsilon_i$  in the physical layer and retransmission limit  $\zeta_i^t \in \mathbb{N}$  in the MAC layer, i.e.  $b_i^t = (\gamma_i^t, \zeta_i^t)$ . Here,  $\Upsilon_i$  is the set of possible modulation schemes. For more sophisticated examples of actions, see e.g. [8] for application layer actions.

### C. State transition and stage reward

Since the network resource state is not affected by the actions performed by the WSTAs, the transition of  $w^t$  can be modeled as a finite state Markov chain (FSMC) [26]. The transition probability is denoted by  $q(w^{t+1} | w^t)$ . In this section, we assume that the transition probability  $q(w^{t+1} | w^t)$  is known by all the WSTAs and CSM. However, more complicated models for the network resource state transition [24] can also be involved in our stochastic game framework.

Deleted: generalized

When WSTA  $i$  receives the resource allocation  $z_i^t$ , it deploys the internal action  $b_i^t$  and can transmit  $n_i^t$  packets during time slot  $t$ , which is computed as

$$n_i^t = \left\lfloor \frac{\Psi(c_i^t, z_i^t, \gamma_i^t, \zeta_i^t) \Delta T}{\ell} \right\rfloor, \quad (6)$$

where  $\Psi(\cdot)$  is the effective rate function, the form of which depends on the protocols implemented at the WSTA. Then, the buffer state can be updated as



$$\begin{bmatrix} v_{i1}^{t+1} \\ \vdots \\ v_{ij}^{t+1} \\ \vdots \\ v_{iJ}^{t+1} \end{bmatrix} = \begin{bmatrix} v_{i2}^t - \max(n_i^t - v_{i1}^t, 0) \\ \vdots \\ v_{i(j+1)}^t - \max\left(n_i^t - \sum_{m=1}^{j-1} v_{im}^t, 0\right) \\ \vdots \\ Y_i^t \end{bmatrix}, \quad (7)$$

where  $Y_i^t$  is a random variable representing the number of packets arriving at time slot  $t$  having life time  $J$ . The distribution of  $Y_i^t$  is denoted by  $p_{Y_i^t}(l)$ . Hence, the transition probability is given by

$$p(v_i^{t+1} | v_i^t, z_i^t, b_i^t) = \begin{cases} P_{Y_i^t}(l) & \text{if } v_i^{t+1} \text{ satisfies eq. (7) and } Y_i^t = l \\ 0 & \text{o.w.} \end{cases} \quad (8)$$

The channel condition  $c_i^t$  depends on the channel gain and the power level for transmission. The channel gain is generally modeled as a FSMC. In this example, we also consider that the power allocation is constant during the data transmission, and hence, the channel condition  $c_i^t$  can be formulated as a FSMC with transition probability  $p(c_i^{t+1} | c_i^t)$ . Details about such transition probability computations can be found in [8].

The state transition probability for WSTA  $i$  is given by

$$p(s_i^{t+1} | s_i^t, z_i^t, b_i^t) = p(v_i^{t+1} | v_i^t, z_i^t, b_i^t) p(c_i^{t+1} | c_i^t). \quad (9)$$

Here, we assume that the transition of the channel condition is independent of the transition of the buffer state. The utility for the delay-sensitive application at time slot  $t$  is defined here as

$$u(s_i^t, z_i^t, b_i^t) = \min\left(n_i^t, \sum_{j=1}^J v_{ij}^t\right) - \lambda_g \min\{v_{i1}^t - n_i^t, 0\}, \quad (10)$$

where  $\lambda_g$  is the parameter to trade-off the received and lost packets (see [22] for details). More sophisticated utility formulations for multimedia transmission, which consider the explicit impact on the multimedia quality (e.g. PSNR) can be found in [22].

#### D. Resource allocation rule

We model the multi-user wireless resource allocation as an auction [7][9][21] for spectrum opportunities held by the CSM during each time slot. The WSTAs calculate the external action  $a_i^t$  based on the information about the network resources, and their own private information about the environment they experience, and their anticipated internal actions [22]. In this auction game, the external action is the competition bid, i.e.  $m_i^t = a_i^t$ . Next, we use the terms - external action and bid interchangeably. Subsequently, each WSTA submits the bid  $a_i^t$  to the CSM. After receiving the bid vectors from the

WSTAs, the CSM computes the channel allocation  $z_i^t$  for each WSTA  $i$  based on the submitted bids. To compel the WSTAs to declare their bids truthfully [25], the CSM also computes the payment  $\tau_i^t \in \mathbb{R}_-$  that the WSTAs have to pay for the use of resources during the current stage of the game. The negative value of the payment represents the absolute value that WSTA  $i$  has to pay the CSM for the used resources. The auction result is then transmitted back to the WSTAs which can deploy their transmission strategies in different layers and send data over the assigned channel. After the data transmission, another auction starts at the next time slot  $t + 1$ . A schematic of the evolution of the multi-user interaction is portrayed for illustration in Figure 6.

The computation of the allocation  $z_i^t$  and payment  $\tau_i^t$  is described as follows. After each WSTA submits the bid vector, the CSM performs two computations: (i) channel allocation and (ii) payment computation. During the first phase, the CSM allocates the resources to WSTAs based on its adopted fairness rule, e.g. maximizing the total “social welfare”<sup>10</sup>:

$$z^{t,opt} = \arg \max_{z^t} \sum_{i=1}^M \tilde{h}_i(a_i^t, z_i^t, w), \quad (11)$$

where  $\tilde{h}_i(\cdot)$  is the utility function of WSTA  $i$  seen by the CSM. Note that this utility can be represented by either the effective rate or time on the network allocated to each user, or it can be determined in the utility domain, by considering the utility-rate functions of the deployed multimedia coders [31].

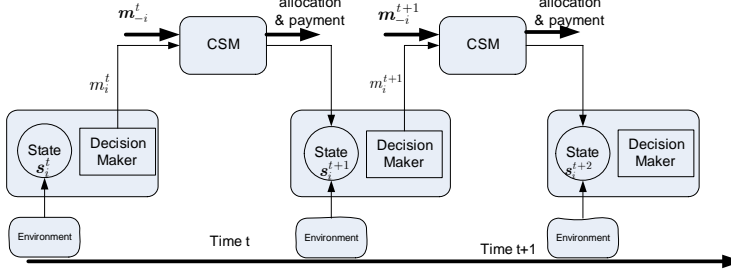


Figure 6. Evolution of multi-user interaction.

We will consider in this paper, for illustration, a second price auction mechanism [20][25] for determining the tax that needs to be paid by WSTA  $i$  based on the above optimal channel assignment  $z^{t,opt}$ . This tax equals:

<sup>10</sup> Note that other social welfare solutions [46] could be adopted and this will not influence our proposed solution.

$$\tau_i^t = \sum_{\substack{j=1, \\ j \neq i}}^M \tilde{h}_j(a_j^t, z_j^{t, opt}, w) - \max_{\substack{z_i^t \\ j \neq i}} \sum_{j=1, \\ j \neq i}^M \tilde{h}_j(a_i^t, z_i^t, w). \quad (12)$$

For simplicity, we can denote the output of the resource allocation game as  $\mathbf{r}^t = (z^t, \tau^t) = \Omega(a^t, w^t)$ .

Note that as mentioned in Section 2, the CSM can design different resource allocation games using different mechanisms, leading to different social decisions, allocations, equilibriums etc. Moreover, the taxation does not need to be implemented and can be omitted. However, we would like to point out that, unless the taxation is implemented, the WSTA will have no incentive to efficiently optimize their cross-layer strategies, upgrade their systems or truthfully and optimally declare their requirements.

#### E. Selecting the Policy for Playing the Resource Management Game

In the cognitive radio network, we assume that the stochastic game is played by all WSTAs for an infinite number of stages. This assumption is reasonable for multimedia applications, which usually have a long duration. In our network setting, we define a history of the stochastic game up to time  $t$  as  $\mathbf{h}^t = \{s^0, w^0, \mathbf{a}^0, \mathbf{b}^0, \mathbf{z}^0, \boldsymbol{\tau}^0, \dots, s^{t-1}, w^{t-1}, \mathbf{a}^{t-1}, \mathbf{b}^{t-1}, \mathbf{z}^{t-1}, \boldsymbol{\tau}^{t-1}, s^t\} \in \mathcal{H}^t$ , which summarizes all previous states and the actions taken by the WSTAs as well as the outcomes at each stage of the auction game and  $\mathcal{H}^t$  is the set of the entire history up to time  $t$ . However, during the stochastic game, each WSTA  $i$  cannot observe the entire history, but rather part of the history  $\mathbf{h}^t$ . The observation of WSTA  $i$  is denoted as  $\mathbf{o}_i^t \in \mathcal{O}_i^t$  and  $\mathbf{o}_i^t \subset \mathbf{h}^t$ . Note that the current state  $s_i^t$  can be always observed, i.e.  $s_i^t \in \mathbf{o}_i^t$ . Then, a bidding policy  $\pi_i^t : \mathcal{O}_i^t \mapsto \mathcal{A}_i \times \mathcal{B}_i$  for WSTA  $i$  at the time  $t$  is defined as a mapping from the observations up to the time  $t$  into the specific action, i.e.  $[a_i^t, b_i^t] = \pi_i^t(\mathbf{o}_i^t)$ . Furthermore, a policy profile  $\boldsymbol{\pi}_i$  for WSTA  $i$  aggregates the bidding policies about how to play the game over the entire course of the stochastic game, i.e.  $\boldsymbol{\pi}_i = (\pi_i^0, \dots, \pi_i^t, \dots)$ . The policy profile for all the WSTAs at time slot  $t$  is denoted as  $\boldsymbol{\pi}^t = (\pi_1^t, \dots, \pi_M^t) = (\pi_i^t, \boldsymbol{\pi}_{-i}^t)$ .

The reward for WSTA  $i$  at the time slot  $t$  is  $R_i^t(s_i^t, \mathbf{r}_i^t, b_i^t) = u(s_i^t, z_i^t, b_i^t) + \tau_i^t$ . Since the resource allocation also depends on other SUs' states and external actions as well as the channel state, the reward can be estimated based on the observation  $\mathbf{o}_i^t$ , and thus, the reward used by a WSTA will be  $R_i^t(s_i^t, \mathbf{o}_i^t, b_i^t)$ . The reward  $R_i^k(s_i^k, \mathbf{o}_i^k, b_i^k)$  at stage  $k$  is discounted by factor  $(\alpha_i)^{k-t}$  at time  $t$ . The factor  $\alpha_i$  ( $0 \leq \alpha_i < 1$ ) is the discounted factor determined by a specific application (for instance, for video streaming applications, this factor can be set based on the tolerable delay). The total discounted sum of rewards  $Q_i^t(\pi_i^t, \boldsymbol{\pi}_{-i}^t | s^t, w^t)$  for SU  $i$  can be calculated at time  $t$  from the state profile  $s^t$  as:

Deleted: generalized

$$Q_i^t((\pi_i^t, \pi_{-i}^t) | s^t, w^t) = \sum_{k=t}^{\infty} (\alpha_i)^{k-t} R_i^k(s_i^k, o_i^k, b_i^k), \quad (13)$$

where  $\pi_{-i}^t(s_{-i}^t) = ([a_j^t, b_j^t]_{j \neq i})$ . We assume that the SUs implement the policy  $\pi^t$  in the subsequent time slots. The total discounted sum of rewards in Eq. (13) consists of two parts: (i) the current stage reward and (ii) the expected future reward discounted by  $\alpha_i$ . Note that SU  $i$  cannot independently determine the above value without explicitly knowing the policies and states of other SUs. The SU maximizes the total discounted sum of future rewards in order to select the bidding policy, which explicitly considers the impact of the current bid vector on the expected future rewards.

We define the *best response*  $\beta_i$  for SU  $i$  to other WSTAs' policies  $\pi_{-i}^t$  as

$$\beta_i(\pi_{-i}^t) = \arg \max_{\pi_i} Q_i^t((\pi_i^t, \pi_{-i}^t) | s^t, w^t) \quad (14)$$

The central issue in such stochastic game in cognitive radio networks is how the best response policies can be determined by the SUs. This will be the topic of Section 4.

### 3.4 Specification of distributed stochastic game

An example of a distributed game is the power-control game played by SUs in the interference channels in cognitive radio network. There are  $M$  SUs, each of which comprises one transmitter and one receiver. There are  $N$  channels potentially vacated by the PUs for SUs transmission. At time slot  $t$ , the network resource state is  $w^t = [w_1^t, \dots, w_N^t] \in \{0, 1\}^N$ . The channel gain of SU  $i$  at channel  $j \in \{1, \dots, N\}$  is  $H_{ii}^j$  and the cross channel gain from transmitter  $i$  (belonging to SU  $i$ ) to receiver  $i'$  (belonging to SU  $i'$ ) at channel  $j$  is  $H_{ii'}^j$ . We assume that the (cross) channel gains for all the SUs are constant.

In this game, the state of SU  $i$  is defined as a vector  $s_i^t \in \{0, 1\}^N$  with each element indicating whether SU  $i$  selects that channel (corresponding to 1) or not (corresponding to 0). The external action  $a_i^t$  of SU  $i$  includes two components: channel selection  $\kappa_i^t$  and power allocation  $\varphi_i^t$ , i.e.  $a_i^t = (\kappa_i^t, \varphi_i^t)$ , where  $\kappa_i^t \subseteq \{1, \dots, N\}$ . For each external action  $a_i^t$  of SU  $i$ , there is power constraint imposed on the power allocation, i.e.

$$\sum_{j \in \kappa_i^t} \varphi_{ij}^t \leq P_i \quad (15)$$

In this power-control game, at the beginning of each time slot, the SUs simultaneously choose the channels over which they will transmit delay-sensitive data and allocate the power on the selected channels under the power constraints. In order not to interfere the PUs, the SUs are not allowed to transmit any data over those channels with  $w_j^t = 0$  (i.e. channel  $j$  is occupied by the PUs). For

simplicity, we consider the case that the SUs are free to choose any channels. Hence, the state of SU  $i$  equals the channel selection action, i.e.  $s_i^t = \kappa_i^t$ . The internal actions for the SUs are empty. The effective transmission rate can be computed as

$$T_i^t(a_i^t, e(a_{-i}^t), w^t) = h_i(s_i^t, a_i^t, b_i^t, e(s_{-i}^t, a_{-i}^t), w^t) = h_i(a_i^t, e(a_{-i}^t), w^t) = \sum_{\substack{j \in \kappa_i^t \\ w_j^t = 1}} \frac{1}{2} \log_2 \left( 1 + \frac{H_{ii}^j \varphi_{ij}^t}{N_{0j} + \sum_{i' \in \mathcal{C}_j^t} H_{ii'}^j \varphi_{i'j}^t} \right), \quad (16)$$

where  $\mathcal{C}_j^t$  is the set of SUs who select channel  $j$  in time slot  $t$ ,  $e(a_{-i}^t) = \sum_{i' \in \mathcal{C}_j^t} H_{ii'}^j \varphi_{i'j}^t$  and  $N_{0j}$  represents the noise level in the selected channel  $j$ . In this power-control game, the stage reward function for SU  $i$  can be defined as effective transmission rate per joule, similarly to [49], i.e.

$$R_i^t(a_i^t, e(a_{-i}^t), w^t) = \frac{T_i^t(a_i^t, e(a_{-i}^t), w^t)}{\sum_{j \in \kappa_i^t} \varphi_{ij}^t}. \quad (17)$$

However, such a reward function cannot satisfy the QoS requirements of multimedia applications. Hence, the following stage reward function can be adopted for such applications:

$$R_i^t(a_i^t, e(a_{-i}^t), w^t) = \frac{\Lambda_i \times \{1 - P_i(T_i^t(a_i^t, e(a_{-i}^t), w^t), d_i)\}}{\sum_{j \in \kappa_i^t} \varphi_{ij}^t}, \quad (18)$$

where  $\Lambda_i$  represents the arrival source rate of the applications of SU  $i$  and  $P_i(T_i^t(a_i^t, e(a_{-i}^t), w^t), d_i)$  represents the packet error rate, which is a function of the effective transmission rate  $T_i^t$  (in equation (16)) and the delay deadline  $d_i$  of the applications of SU  $i$ . More complicated utility-resource functions as in [22] can also be employed.

Since the state of each SU  $i$  is the same as the channel selection and no internal actions are considered, the channel-selection and power-control game is reduced to a repeated game [39]. The essential goal of SU  $i$  is to find the best response to the aggregated interference  $e(a_{-i}^t)$  under various network resource states, i.e.

$$a_i^{t,*}(e(a_{-i}^t), w^t) = \arg \max_{a_i^t} R_i^t(a_i^t, e(a_{-i}^t), w^t) \quad (19)$$

This will be discussed in the next section.

## 4. Strategic learning solutions in multi-user wireless systems

### 4.1. Why learn?

In the previous section, it was shown that in order for an SU to derive its own transmission policy, it needs to know how its decision process and resulting performance is coupled to that of other SUs. In a stochastic game framework, the goal for each SU is to find a policy  $\pi_i$  such that its own utility is

Deleted: generalized

maximized. However, as discussed in Section 3, SU  $i$ 's policy  $\pi_i$  depends on other SUs' policies, which is formulated as:

$$\pi_i^* = \arg \max_{\pi_i} Q_i(\pi_i, \boldsymbol{\pi}_{-i} \mid s_i, \mathbf{s}_{-i}, \mathbf{w}), \quad (20)$$

To solve this optimization, the following information is required by SU  $i$ :

- the state transition model of SU  $i$ ,  $p(s_i^{t+1} \mid s_i^t, a_i^t, \mathbf{a}_{-i}^t, b_i)$ ;
- the state transition model of other SUs,  $p(s_j^{t+1} \mid s_j^t, a_j^t, \mathbf{a}_{-j}^t, b_j), \forall j \neq i$ ;
- the state of other SUs,  $\mathbf{s}_{-i}$ ;
- the policy of other SUs,  $\boldsymbol{\pi}_{-i}$ ;
- the network resource state  $\mathbf{w}$ .

This coupling among SUs is due to the shared nature of the wireless resources [2]. However, an SU may not exactly know the other SUs' actions and models, and it cannot know their private information. Thus, an SU can only predict these dynamics (uncertainties) caused by the competing SUs based on its observations from past interaction. In the cognitive radio network, there are different levels of information availability:

- *Private information*: this private information includes the characteristics of the application traffic, channel gain or channel conditions (SINR, etc.).
- *Network information*: the network information refers to the network resource states or the behaviors of PUs.
- *Opponents information*: this information includes the states and possible actions of the opponents.

This information can be for instance known when all the SUs adopt the same protocol, having the same set of states and actions.

To reduce the uncertainty and increase the knowledge about the environment when selecting an action, an SU can deploy learning in games algorithms [23]. Depending on the information availability, different learning solutions can be deployed by a WSTA. The existing learning in games literature provides a broad spectrum of analytical and practical results on learning algorithms and underlying game structures for a variety of competitive interaction scenarios. In general, the main issue considered has been to characterize long term behavior in terms of a generalized equilibrium concept or characterize the *lack* of convergence for general classes of learning dynamics. However, when selecting learning solutions for wireless networks games, the specific constraints and features of wireless systems will need to be considered. For instance, the learning algorithm that should be deployed by a user in a wireless

environment strongly depends on what information an SU can observe about the other SUs, given the adopted protocols or spectrum regulation rules. Moreover, unlike a majority of work in learning in games solutions [23], where the main focus is on proving the existence of equilibriums, or where the only goal of the agents is to achieve different equilibrium conditions, learning solutions in wireless networks are deployed by self-interested and heterogeneous users, which have as only goal to improve their own performance. Thus, **a learning algorithm  $\mathcal{L}_i$  adopted by SU  $i$  to efficiently play the spectrum allocation game will be evaluated based on the information that can be acquired (i.e. the observed information  $o_i^t$ ) and exchanged  $I_{-i}^t$ , the complexity requirements, and the resulting (long-term or short-term) utility  $U_i$ .**

#### 4.2. Definitions of learning algorithms and beliefs

The goal of learning for an SU in the multi-user games is to update its own policy and belief about the other SUs' states and policies. Specifically, by learning from the observed and exchanged information, a user can build its *belief* on the other users' strategies, and determine its own best response policy. In our stochastic game framework, the SU also needs to update its knowledge about the network resource state using learning. We note that a learning algorithm is built based on the observation  $o_i^t$  and exchanged information  $I_{-i}^t$  and hence, it is denoted as  $\mathcal{L}_i(o_i, I_{-i})$  where  $o_i, I_{-i}$  are all the observation and exchanged information obtained by SU  $i$ .

A learning algorithm  $\mathcal{L}_i$  can be defined using the following equations:

$$[a_i^t, b_i^t] = \pi_i^t(s_i^t, B_{s_{-i}}^t, B_{\pi_{-i}}^t, B_w^t) \quad (21)$$

$$\Omega^t = \text{Game}(s^t, a^t, w^t) \quad (22)$$

$$o_i^t = O(s_i^t, \Omega_i^t, b_i^t) \quad (23)$$

$$\pi_i^{t+1} = \mathcal{F}_i(\pi_i^t, o_i^t, I_{-i}^t) \quad (24)$$

$$B_{\pi_{-i}}^{t+1} = \mathcal{F}_{\pi_{-i}}(B_{\pi_{-i}}^t, o_i^t, I_{-i}^t) \quad (25)$$

$$B_w^{t+1} = \mathcal{F}_w(B_w^t, o_i^t, I_{-i}^t) \quad (26)$$

$$B_{s_{-i}}^{t+1} = \mathcal{F}_{s_{-i}}(B_{s_{-i}}^t, o_i^t, I_{-i}^t) \quad (27)$$

where  $B_{s_{-i}}^t$ ,  $B_{\pi_{-i}}^t$  and  $B_w^t$  are the belief about the other SUs' states  $s_{-i}$ , policies  $\pi_{-i}$  and the network resource state  $w$ , respectively;  $\Omega^t$  is the output of the multi-user interaction game ( $\Omega^t = \mathbf{r}^t$  in the centralized stochastic game and  $\Omega^t = \{a_i^t, e(\mathbf{a}_{-i}^t), w^t\}$  for the distributed stochastic game or repeated game);  $o_i^t$  is the observation of SU  $i$  and  $O$  is the observation function which depends on the current

state, the current game output and the current internal action taken;  $\mathcal{F}$  is the update function about the belief and policies;  $I_{-i}^t$  is the exchanged information with the other SUs.

Eq. (21) shows that SU  $i$  generates the external actions based on its own states, the belief about the other SUs' states, policies and network resource state. After each SU executes its external actions, a multi-user spectrum access game is played and the results of the game are produced as shown in Eq. (22). The results of the multi-user game may or may not be fully observed by the SUs based on the game form or the implemented network protocol. Eq. (23) represents the observation function which depends on the network protocols and the SUs' measurement methods. Different (accurate or inaccurate) observations may lead to different learning algorithms, which will be discussed in subsequent sections. Hence, an SU may have incentives to exchange information with other SUs. The exchanged information  $I_{-i}^t$  may be used to update the belief about the other SUs' states, policies as well as the network resource state. Eqs. (24)-(27) represent the updates of the beliefs.

In a wireless communication game, we differentiate two types of users based on their response strategies:

- *Myopic users*: Users that always act to maximize their immediate achievable reward. They are myopic in the sense that, at each decision stage, they treat other users' actions as fixed, ignore the impact of its competitors' reactions over its own performance, and determine their responses to gain the maximal immediate rewards.
- *Foresighted users*: Users that behave by taking into account the long-term impacts of their actions on their rewards. They avoid shortsighted (myopic) actions, anticipate how the other users will react, and maximize their performance by considering the responses of the other users [8][30]. Note that such foresighted users require additional knowledge about the other users to assist their decision making. We will discuss this in more detail later in this section.

Before we proceed in detail with discussing how a learning algorithm is built, we discuss first how we can evaluate a learning algorithm for the cognitive radio network.

#### **4.3. Value of learning, value of information and regret computation**

As mentioned previously, the performance of a learning algorithm will depend on the resulting SU reward. We denote a policy generated by the learning algorithm  $\mathcal{L}_i$  as  $\pi_i^{\mathcal{L}_i}$ . An SU will learn in order to improve its policy and its rewards from participating in the spectrum access game. The performance of SU  $i$  when adopting the learning algorithm  $\mathcal{L}_i$  is defined as the time average reward obtained in a time



window with length  $T$  in which this learning algorithm was used:

$$\mathcal{V}_{\pi_i}^{\mathcal{L}_i(o_i, I_{-i})}(T) = \frac{1}{T} \sum_{t=1}^T R_i^t(\pi_i^{\mathcal{L}_i(o_i, I_{-i})}), \quad (28)$$

where the reward  $R_i^t$  depends on both the learning approach  $\mathcal{L}_i$  and on the observation  $o_i^t$  and information exchanged  $I_{-i}^t$ . Thus, using this definition, the “value of a learning scheme” can be determined. For instance, given the same observation  $o_i^t$  and exchanged information  $I_{-i}^t$ , if the time average rewards of two algorithms  $\mathcal{L}_i'$  and  $\mathcal{L}_i''$  satisfy  $\mathcal{V}_{\pi_i}^{\mathcal{L}_i'(o_i, I_{-i})}(T) > \mathcal{V}_{\pi_i}^{\mathcal{L}_i''(o_i, I_{-i})}(T)$ , then we say that learning algorithm  $\mathcal{L}_i'$  is better than  $\mathcal{L}_i''$ . The “value of information exchange” with respect to a learning algorithm  $\mathcal{L}_i$  can be also similarly computed. This value of information will play a significant role on what information should be exchanged among WSTAs and how WSTAs should negotiate in a cooperative<sup>11</sup> setting (e.g. in a bargaining or coalition setting). The value of making various observations and learning based on them can be similarly computed. Moreover, similar to [47], we define a generalized “regret” for the stochastic game at each time slot  $t$  as

$$\Delta_{\mathcal{L}_i} \triangleq \max_{\pi_i} Q_i^t(\pi_i, \pi_{-i} \mid s_i^t, s_{-i}^t, w^t) - Q_i^t(\pi_i^{\mathcal{L}_i}, B_{\pi_{-i}}^t \mid s_i^t, B_{s_{-i}}^t, B_w^t). \quad (29)$$

When the stochastic game is reduced to a repeated game, the regret can be computed as

$$\Delta_{\mathcal{L}_i} \triangleq \max_{a_i^t} R_i^t(a_i^t, e(a_{-i}^t), w^t) - R_i^t(a_i^{\mathcal{L}_i}, B_{a_{-i}}^t, B_w^t) \quad (30)$$

The regret is computed as the reward loss due to the lack of knowledge about the network resource and components’ states and actions. The regret can be computed and used by the SUs in order to adjust their learning strategies and improve their strategies for playing the game.

#### 4.4. Learning framework for wireless stochastic games

One possible simplification for the stochastic learning is to assume that other SUs perform a fixed policy. This is a good assumption especially for the case when WSTAs adopt the same protocols, which implement the same policies. Hence, SU  $i$  does not need to update its belief about other SUs’ policies (i.e.  $B_{\pi_{-i}}^t$ ). Instead, SU  $i$  needs to update its belief about other SUs’ states and state transition probability. However, to observe other SUs’ states in cognitive radio network is also difficult, and even impossible in some cases. To solve this problem, an SU can classify the states of other SUs based on the output of the game. **For simplification, we assume that the network resource state is common information known by all**

Formatted: Font: Not Italic

<sup>11</sup> The term “cooperative” (in the cooperative game theory) does *not* mean that decision makers have interests that are completely aligned. Rather, cooperative game solution concepts are relevant in situations where a scarce resource is to be rationed *fairly* among competing claimants that are strategically negotiating with each other, as in bargaining and coalition games.

participating SUs. However, the learning algorithm discussed in this section can also be extended to the case in which the network resource state and the corresponding state transition are unknown to the SUs. In this case, the WSTA needs to learn the state transition probability for each channel's state based on the observation [51].

Formatted: Font: Not Italic

Formatted: Font: Not Italic

**Deleted:** For simplification, we assume that the network resource state is common information known by all participating SUs. However, the learning algorithm discussed in this section can also be extended to the case in which the network resource state is unknown for the SUs.

#### A. What Information to Learn from?

First, let us consider what information the SU can observe while playing the stochastic game in our cognitive radio network. As shown in Figure 3, at the beginning of time slot  $t$ , the SUs submit the bids  $a_i^t, \forall i$ . Then, the CSM returns the channel allocation  $z_i^t, \forall i$  and  $\tau_i^t, \forall i$ . In cognitive radio network, if SU  $i$  is not allowed to observe the bids, the channel allocations and payments for other SUs, then the observation of SU  $i$  becomes  $o_i^t = \{s_i^0, w^0, a_i^0, b_i^0, z_i^0, \tau_i^0, \dots, s_i^{t-1}, w^{t-1}, a_i^{t-1}, b_i^{t-1}, z_i^{t-1}, \tau_i^{t-1}, s_i^t\}$ . If the information is fully exchanged among SUs or broadcasted and overheard by all SUs, the observed information by SU  $i$  becomes  $o_i^t = h^t$ . Now, the problem that needs to be solved by SU  $i$  is how it can improve its own policy for playing the game by learning from the observation  $o_i^t$ . In this paper, we assume that SU  $i$  observes the information  $o_i^t = \{s_i^0, w^0, a_i^0, b_i^0, z_i^0, \tau_i^0, \dots, s_i^{t-1}, w^{t-1}, a_i^{t-1}, b_i^{t-1}, z_i^{t-1}, \tau_i^{t-1}, s_i^t\}$ .

#### B. What needs to be learnt?

A key question is what needs to be learned within a wireless stochastic game in order to improve the policy of an SU. We focus here on the learning procedure for the external policy (generating external actions, i.e. bidding actions).

In Section 4.1 we discussed the information that SU  $i$  needs to learn in order to be able to solve the optimization in eq. (20). However, SU  $i$  can only observe the information  $o_i^t = \{s_i^0, w^0, a_i^0, b_i^0, z_i^0, \tau_i^0, \dots, s_i^{t-1}, w^{t-1}, a_i^{t-1}, b_i^{t-1}, z_i^{t-1}, \tau_i^{t-1}, s_i^t\}$  from which SU  $i$  cannot accurately infer the other SUs' state space (i.e.  $\mathcal{S}_{-i}$ ), the current state of other SUs (i.e.  $s_{-i}^t$ ) and the transition probability of other SUs (i.e.  $\prod_{k \neq i} q_k(s_k^{t+1} | s_k^t, z_k^t)$ ). Moreover, capturing the exact information about other SUs requires heavy computational and storage complexity. Instead, we allow SU  $i$  to classify the space  $\mathcal{S}_{-i}$  into  $H_i$  classes, each of which is represented by a representative state  $\tilde{s}_{-i,h}, h \in \{1, \dots, H_i\}$ . By dividing the state space  $\mathcal{S}_{-i}$ , the transition probability  $\prod_{k \neq i} q_k(s_k^{t+1} | s_k^t, z_k^t)$  is approximated by  $q_{-i}(\tilde{s}_{-i}^{t+1} | \tilde{s}_{-i}^t, z_i^t)$ , where  $\tilde{s}_{-i}^t$  and  $\tilde{s}_{-i}^{t+1}$  are the representative states of the classes that  $s_{-i}^t$  and  $s_{-i}^{t+1}$  belong to. This approximation is performed by aggregating all other SUs' states into one representative state and assuming that the transition depends on the resource allocation  $z_i^t$ . Note that the classification on the state space  $\mathcal{S}_{-i}$  and approximation of the transition probability and discounted sum of rewards affects the

learning performance. Hence, a user should tradeoff an increased learning complexity for an increased value of learning.

In this setting, to find the approximated optimal bidding policy, we need to learn the following from the past observations: (i) how the space  $\mathcal{S}_{-i}$  is classified; (ii) the transition probability  $q_{-i}(\tilde{s}_{-i}^{t+1} | \tilde{s}_{-i}^t, \mathbf{z}_i^t)$ ; (iii) the average future rewards  $V_i^{t+1}((\mathbf{s}_i^{t+1}, \tilde{s}_{-i}^{t+1}))$ .

### C. How to learn?

In this section, we develop a learning algorithm to estimate the terms listed in the above section.

#### Step 1. Decomposition of the space $\mathcal{S}_{-i}$

As discussed in Section A, only  $\mathbf{o}_i^t = \{\mathbf{s}_i^0, w^0, \mathbf{a}_i^0, \mathbf{b}_i^0, \mathbf{z}_i^0, \boldsymbol{\tau}_i^0, \dots, \mathbf{s}_i^{t-1}, w^{t-1}, \mathbf{a}_i^{t-1}, \mathbf{b}_i^{t-1}, \mathbf{z}_i^{t-1}, \boldsymbol{\tau}_i^{t-1}, \mathbf{s}_i^t\}$  are observed. From the auction mechanism presented in Section 3.D, we know that the value of the tax  $\tau_i^t$  is computed based on the inconvenience that SU  $i$  causes to the other SUs. In other words, a higher value of  $|\tau_i^t|$  indicates that the network is more congested<sup>12</sup>. Based on the bid vector  $\mathbf{b}_i^t$ , the channel allocation  $\mathbf{z}_i^t$  and the tax  $\tau_i^t$ , SU  $i$  can infer the network congestion and thus, indirectly, the resource requirements of the competing SUs. Instead of knowing the exact state space of other SUs, SU  $i$  can classify the space  $\mathcal{S}_{-i}$  as follows. We assume the maximum absolute tax is  $\Gamma$ . We split the range  $[0, \Gamma]$  into  $[\Gamma_0, \Gamma_1), [\Gamma_1, \Gamma_2), \dots, [\Gamma_{H_i-1}, \Gamma_{H_i}]$  with  $0 = \Gamma_0 \leq \Gamma_1 \leq \dots \leq \Gamma_{H_i} = \Gamma$ . Here, we assume that the values of  $\{\Gamma_1, \dots, \Gamma_{H_i-1}\}$  are equally located in the range of  $[0, \Gamma]$ . (Note that more sophisticated selection for these values can be deployed, and this forms an interesting area of future research.)

We need to consider three cases to determine the representative state  $\tilde{s}_{-i}^t$  at time  $t$ .

(1) If the resource allocation  $\mathbf{z}_i^t \neq \mathbf{0}$ , then the representative state of other SUs is chosen as

$$\tilde{s}_{-i}^t = h, \text{ if } |\tau_i^t| \in [\Gamma_{h-1}, \Gamma_h). \quad (31)$$

(2) If the resource allocation  $\mathbf{z}_i^t = \mathbf{0}$  but  $w^t \neq 0$ , the tax is 0. In this case, we cannot use the tax to predict the network congestion. However, we can infer that the congestion is more severe than the minimum bid for those available channels, i.e.  $\min_{j \in \{l: y_l^t \neq 0\}} \{b_{ij}^t\}$ . This is because, in this current stage of the auction game, only SU  $i'$  with  $b_{i'j}^t \geq b_{ij}^t$  can obtain channel  $j$  which indicates that  $|\tau_i^t| \geq \min_{j \in \{l: y_l^t \neq 0\}} \{b_{ij}^t\}$ , if SU  $i$  is allocated any channel. Then the representative state of other SUs is chosen as

$$\tilde{s}_{-i}^t = h, \text{ if } \min_{j \in \{l: y_l^t \neq 0\}} \{b_{ij}^t\} \in [\Gamma_{h-1}, \Gamma_h) \quad (32)$$

(3) If the resource allocation  $\mathbf{z}_i^t = \mathbf{0}$  and  $w^t = 0$ , there is no interaction among the SUs in this time slot. Hence,  $\tilde{s}_{-i}^t = \tilde{s}_{-i}^{t-1}$ .

<sup>12</sup> When the CSM deploys a mechanism without tax for the resource management, the space classification for other SUs can also be done based on the announced information and corresponding resource allocation.

### Step 2. Estimating the transition probability

To estimate the transition probability, SU  $i$  maintains a table  $F$  with size  $H_i \times H_i \times (N+1)$ . Each entry  $f_{h',h'',j}$  in the table  $F$  represents the number of transitions from state  $\tilde{s}_{-i}^t = h''$  to  $\tilde{s}_{-i}^{t+1} = h'$  when the resource allocation  $\mathbf{z}_i^t = \mathbf{e}_j$  (or  $\mathbf{0}$  if  $j = 0$ ). Here,  $\mathbf{e}_j$  is a  $N$ -dimensional vector with the  $j$ -th element being 1 and otherwise being 0. It is clear that  $H_i$  will influence significantly the complexity and memory requirements etc. of SU  $i$ . The update of  $F$  is simply based on the observation  $\mathbf{o}_i^t$  and the state classification in the above section. Then, we use the frequency to approximate the transition probability [17], i.e.

$$q_{-i}(\tilde{s}_{-i}^{t+1} = h' \mid \tilde{s}_{-i}^t = h'', \mathbf{e}_j) = \frac{f_{h',h'',j}}{\sum_{h'} f_{h',h'',j}} \quad (33)$$

### Step 3. Learning the future reward

By classifying the state space  $\mathcal{S}_i$  and estimating the transition probability, SU  $i$  can now forecast the value of the average future reward  $V_i^{t+1}((\mathbf{s}_i^{t+1}, \tilde{s}_{-i}^{t+1}))$  using learning. Eq. (13) can be approximated by

$$Q_i^t(\mathbf{s}_i^t, \tilde{s}_{-i}^t) = \{g_i(\mathbf{s}_i^t, \mathbf{z}_i^t) + \tau_i^t + \alpha_i \sum_{(\mathbf{s}_i^{t+1}, \tilde{s}_{-i}^{t+1}) \in \mathcal{S}} q_i(\mathbf{s}_i^{t+1} \mid \mathbf{s}_i^t, \mathbf{z}_i^t) q_{-i}(\tilde{s}_{-i}^{t+1} \mid \tilde{s}_{-i}^t, \mathbf{z}_i^t) V_i^{t+1}((\mathbf{s}_i^{t+1}, \tilde{s}_{-i}^{t+1}))\} \quad (34)$$

The received rewards are used to update the estimation of future rewards, similarly to Q-learning [19]. However, the main difference between this algorithm [8] and Q-learning is that the former explicitly considers the impact of other SUs' bidding actions through the state classifications and transition probability approximation.

A 2-dimensional table can be used to store the value  $V_i(\mathbf{s}_i, \tilde{s}_{-i})$  with  $\mathbf{s}_i \in \mathcal{S}_i$ ,  $\tilde{s}_{-i} \in \tilde{\mathcal{S}}_{-i}$ , where  $\tilde{\mathcal{S}}_{-i}$  is the set of representative states for the other SUs. The total number of entries in  $V_i$  is  $|\mathcal{S}_i| \times |\tilde{\mathcal{S}}_{-i}|$ . SU  $i$  updates the value of  $V_i((\mathbf{s}_i, \tilde{s}_{-i}))$  at time  $t$  according to the following rules:

$$V_i^t(\mathbf{s}_i, \tilde{s}_{-i}) = \begin{cases} (1 - \gamma_i^t) V_i^{t-1}(\mathbf{s}_i, \tilde{s}_{-i}) + \gamma_i^t Q_i^t(\mathbf{s}_i^t, \tilde{s}_{-i}^t) & \text{if } (\mathbf{s}_i^t, \tilde{s}_{-i}^t) = (\mathbf{s}_i, \tilde{s}_{-i}) \\ V_i^t(\mathbf{s}_i, \tilde{s}_{-i}) & \text{otherwise} \end{cases} \quad (35)$$

where  $\gamma_i^t \in [0,1)$  is the learning rate factor. An interesting area of research is determining how the learning rate factor should be determined (and possibly adapted) in various cognitive radio settings, where different dynamics are experienced.

### D. Complexity of the learning algorithm

In this section, we quantify the complexity of learning in terms of the computational and storage burden. We use the ‘‘flop’’ (floating-point operation) as a measure of complexity, which will provide us an estimation of the computational complexity required for performing the learning algorithm. Also, based on this, we can determine how the complexity grows with the increasing number of SUs. At each stage,

Formatted: Font: Times New

Formatted: Heading 2, H2, H21, © o, µ 2, Left, Don't adjust space between Latin and Asian text, Don't adjust space between Asian text and numbers

Formatted: Bullets and Numbering

the SU performs the classification of other SUs' states, which, in the worst case, requires a number of "flops" of approximately  $N$ . The number of "flops" for estimating the transition probability of other SUs' states as in Eq. **Error! Reference source not found.** is approximately  $(H_i + 1)$ . The number of "flops" for learning the future reward is approximately  $(2|S_i|H_i + 6)$ . Therefore, the total number of "flops" incurred by the SU is  $N + H_i + 2|S_i|H_i + 7$ , from which we can note that the complexity of learning for each SU is proportional to the number of possible states of that SU and the number of classes in which the other SUs' state space is decomposed.

To perform the learning algorithm, the SU needs to store 2 tables (i.e. transition probability table and state-value table) which have totally  $(H_i^2(N + 1) + 2^N|S_i|H_i)$  entries. We also note that the storage complexity is also proportional to the number of possible states of that SU and the number of classes in which the other SUs' state space is decomposed.

**Formatted:** Style Style Text + First line: 0.17" Char Char + 11 pt, Left, Line spacing: single, Widow/Orphan control

#### 4.5. Illustration of various bidding and learning strategies

In this section, we highlight the performance of the learning framework presented in the previous section in a centralized stochastic game (introduced in Section 3). We assume that the SUs compete for the available spectrum opportunities in order to transmit delay-sensitive multimedia data. The SUs can deploy different bidding strategies to generate their bid vector:

- Fixed bidding strategy  $\pi_i^{fixed}$ : this strategy generates a constant bid vector during each stage of the auction game, irrespective of the state that SU  $i$  is currently in and of the states other SUs are in. In other words,  $\pi_i^{fixed}$  does not consider any source and channel dynamics.
- Source-aware bidding strategy  $\pi_i^{source}$ : this strategy generates various bid vectors by considering the dynamics in source characteristics (based on the current buffer state), but not the channel dynamics.
- Myopic bidding strategy  $\pi_i^{myopic}$ : this strategy takes into account both the environmental disturbances and the impact caused by other SUs. However, it does not consider the impact on its future rewards.
- Bidding strategy based on best response learning  $\pi_i^{\mathcal{L}}$ : This strategy is produced using the learning algorithm presented in the previous section, which considers both the environmental dynamics and the interaction impact on the future reward.

In this simulation, we consider the cognitive radio network as an extension of WLANs with cognitive radio capability [11]. (More simulation details can be found in [8].) To highlight the impact on the multimedia quality, in this illustrative simulation, we assume that both users are streaming to their receivers the *Coastguard* video sequence and they both tolerate an application layer delay of 500ms. For

illustration, the following four scenarios are considered. In scenario (1)~(4), SU 1 deploys a fixed bidding strategy  $\pi_1^{fixed}$ , source-aware bidding strategy  $\pi_1^{source}$ , myopic bidding strategy  $\pi_1^{myopic}$  and best response learning based bidding strategy  $\pi_1^{\mathcal{L}}$ , respectively, and SU 2 always deploys the myopic bidding strategy  $\pi_2^{myopic}$ . The average video quality (PSNR), average tax and average reward per time slot (see Section 3.E) are presented in Table 1.

From this simulation, we observe that when SU 2 deploys the myopic strategy, SU 1 increase its own reward by adopting advanced learning algorithms (from fixed bidding strategy  $\pi_i^{fixed}$  to best response learning based bidding strategy  $\pi_i^{\mathcal{L}}$ ). On the other hand, SU 2 starts to have an increased cost as SU1 starts to deploy increasingly advanced learning algorithms.

It is also worth to note that the improvement in video quality for SU 1 in scenarios 1~4 comes from two parts: one is the advanced bidding strategies, which allows the SU to take into consideration more information about its own states and the other SUs' states and, based on this, better forecast the impact of various actions, and the other one is the increase in the amount of resources consumed by SU 1 which corresponds to higher tax charged by the CSM, as shown in Table 1.

Table 1. Performance of SU 1 and 2 with various bidding strategies by the two competing SUs

	Bidding Strategies	SU 1			SU 2		
		Video Quality (PSNR)	Average tax	Average reward	Video Quality (PSNR)	Average tax	Average reward
Scenario 1	$\pi_1^{fixed}, \pi_2^{myopic}$	25 dB	0.1222	2.6337	36 dB	0.5495	1.5105
Scenario 2	$\pi_1^{source}, \pi_2^{myopic}$	26 dB	0.3147	2.4915	33 dB	0.6048	1.6116
Scenario 3	$\pi_1^{myopic}, \pi_2^{myopic}$	29 dB	0.4669	1.9767	30 dB	0.3763	1.7837
Scenario 4	$\pi_1^{\mathcal{L}}, \pi_2^{myopic}$	35 dB	0.6923	1.7428	27 dB	0.4197	2.2967

#### 4.6 Learning in repeated games

A simplification of the stochastic game is the case where each SU has only one state. In this case, the stochastic game is reduced to a repeated game. In this case, the policy for each SU becomes the same as the action that each SU selected. Thus, an SU only needs to update its belief about the other SUs' actions.

##### A. Myopic adaptation

In wireless communication, a simple learning (or adaptation) method is myopic adaptation, where the SU does not update its belief about the other SUs' actions. Instead, it maximizes its utility based on the aggregated observation of other SUs' actions during the previous round of game, i.e.

$$a_i^{t,myopic} = \arg \max_{a_i^t} R_i^t(a_i^t, O(a_{-i}^{t-1}), w^t), \quad (36)$$

where  $O(a_{-i}^{t-1})$  represents the aggregated observation of other SUs' actions in time slot  $t-1$ .

In power control games among WSTAs in interference channels, the myopic adaptation has been proved to converge to the Nash equilibrium point [49], which generally leads to a lower system performance for the user than the collaborative case, where a moderator will compel the WSTAs to operate on the Pareto surface.

## B. Reinforcement learning

Comment [F1]: I don't think it is a good title.

In the reinforcement learning solution, an SU does not need to know the actions of the other SUs. Hence, this method is very suitable in a variety of repeated wireless games, including the abovementioned power control games [38]. In this learning, the SU establishes a preference for each action. The preference is updated based on the utility that it obtains during the different stages of the game, without trying to explicitly model the other SUs' actions. Then, based on its preference, the SU determines a mixed action to perform during each time slot. Formally, when adopting the reinforcement learning algorithm, SU  $i$  computes its best response mixed action  $A_i^t$  as

$$A_i^t(a_i) = \frac{\phi(\rho_i^t(a_i))}{\sum_{a_i \in \mathcal{A}} \phi(\rho_i^t(a_i))}, \quad (37)$$

where  $\rho_i^t(a_i)$  represents the *preference* [47] of SU  $i$  choosing an action  $a_i$  at time slot  $t$ , and  $\phi(\cdot)$  is a non-decreasing positive function (e.g.  $\phi(x) = e^x$ ), and  $A_i^t(a_i)$  is the mixed action. When an action  $a_i$  is adopted by SU  $i$  at time slot  $t$ , the reward  $R_i^t(a_i, \mathbf{a}_{-i})$  is obtained. This reward is used to update the preference as follows:

$$\rho_i^t(a_i) = \rho_i^{t-1}(a_i) + \alpha [R_i^t(a_i, \mathbf{a}_{-i}) - \rho_i^{t-1}(a_i)], \quad (38)$$

where  $\alpha$  is a update step size. An adaptive reinforcement (AR) technique can also be implemented, in which an SU can adapt its preference with various frequencies corresponding to different learning speeds, based on a cost-benefit tradeoff. A faster learning speed provides more accurate belief updates (in equations (25)~(27)), however it also requires a slightly higher computational cost and higher private information feedback overheads associated with the increased observations (in equation (23)).

## C. Action-based learning

In this setting, an SU explicitly models the exact actions of other SUs by directly exchanging information with other SUs (i.e.  $I_{-i}^t$ ) about their taken actions. In this case, fictitious play and regret matching solutions can be used [47]. For instance, SU  $i$  can adopt an adaptive fictitious play algorithm, where it maintains a set of strategy vectors  $\mathbf{a}_{-i}^t[a_{-i} | a_i] = \{\mathbf{a}_j^t[a_j \in \mathcal{A}_j | a_i], \text{ for all SUs } j \neq i\}$  for all

possible actions  $a_i \in \mathcal{A}_i$ , with  $\mathbf{a}_j^t[a_j \in \mathcal{A}_j \mid a_i]$  representing the estimated strategy of the other users  $j \neq i$  given that SU  $i$  took action  $a_i$  at time slot  $t$ . The adaptive fictitious play algorithm models the actions of other SUs  $j \neq i$  as:

$$A_j^t(a_j \mid a_i) = \frac{\phi(\rho_j^t(a_j \mid a_i))}{\sum_{a_j \in \mathcal{A}_j} \phi(\rho_j^t(a_j \mid a_i))}, \quad (39)$$

where  $\rho_j^t(a_j \mid a_i)$  represents the *anticipating preference* of SU  $j$  choosing an action  $a_j$  at time slot  $t$ , given that the anticipator SU  $i$  taking an action  $a_i$ . The preference can be updated similarly as in the reinforcement learning case. Moreover, adaptive versions of this action learning, which we refer to as adaptive action (AA) learning, can also be adopted, where an SU is modeling other SUs with different accuracies in order to reduce the informational overhead and the computational overhead. This is especially important in the dynamic power/spectrum management games, where the neighboring SUs can be classified by an SU based on their impact on its utility. For instance, a neighboring SU with a larger channel gain will have higher impact on its utility.

#### 4.7. Illustrative results for different learning approaches in repeated games

Next, we show several illustrative results using the learning schemes discussed in the previous sections in the distributed power control repeated games. We assume that 5 SUs (distinct transmitter-receiver pairs) are in the network and share 3 frequency channels. Each user can choose its power level from a set  $\mathcal{P} = \{20, 40, 60, 80, 100\}$  (mW). Hence, there are a total of 15 actions for users to select. For the application layer parameters, we set the average packet length  $L_v = 1000$  bytes, input rate  $R_v = 500$  Kbps ( $\Lambda_v = R_v / L_v$ ), and delay deadline  $d_v = 200$  msec for all the users.

Besides the Adaptive Reinforcement (AR) scheme mentioned in Section 4.6B and the Adaptive Action learning (AA) scheme mentioned in Section 4.6C, we also consider the myopic best response without learning discussed in Section 4.6A, which leads to a Nash Equilibrium (NE). We select SU 1 to be the user who learns from the observed information.

The results are presented in Table 2, where the reward is defined as in equation (18). From the results, it is interesting to see how the resulting reward of SU 1 improves when this user starts learning (when using AA and AR scheme) as opposed to the case that it is merely adopting a myopic best response (when using NE scheme). Using the AA scheme, users are able to exploit the spectrum more efficiently, due to the ability that the users can better model the strategies of other interference sources in the network.



However, this requires significant information overhead, which results in a worse performance than using the AR scheme. Note that although only SU 1 is learning, the average reward of using interactive learning schemes outperforms the myopic NE scheme. Thus, as discovered in [30], this foresighted user benefits both itself as well as the overall system performance.

Table 2. Simulation results for various repeated games, using different learning techniques.

Adopted schemes	SU	Reward (Kbit/joule)	Average reward
Myopic scheme	1	519.0	890.15
	2	195.2	
	3	530.6	
	4	2073.0	
	5	1132.9	
AR learning scheme	1	555.2	1005.6
	2	113.5	
	3	345.6	
	4	2830.2	
	5	1183.7	
AA learning scheme	1	529.3	1069.3
	2	475.6	
	3	476.8	
	4	2831.2	
	5	1033.3	

#### 4.8. Future research directions for learning in communication networks

Learning in games offers significant potential as a paradigm for shaping dynamic wireless network interactions. As stated earlier, the majority of the research literature in this topic was aimed at proving that different types of equilibriums exist [23]. However, in wireless networks the focus is on constructing adaptive algorithms and protocols that allow SUs to interact with each other based on their knowledge level in order to improve their performance. Accordingly, there are important research directions remaining to be addressed to enable the SUs and the wireless system to achieve the optimal performance. In particular, typical assumptions on knowledge of utility functions in multi-agent learning are of the “all or nothing” type. That is, either an agent knows the utility function fully or can only measure payoffs online. A middle ground is the case where there is partial knowledge of the functional form, but subject to uncertain parameters that may be estimated online. For instance, in the discussed communication setting, users sharing the same protocol have the same states and actions. The only difference is that they experience different private information. Thus, model-based learning approaches can be deployed that take advantage of the fact that users in the same protocol class adopt the same utility functions. These methods allow a user to learn more effectively, since they only need to learn the model parameters. Moreover, this can be also extended to the case where both the parameters and the models are unknown.

## 5. Conclusions

In this paper, we provided a unifying framework for dynamic spectrum access and learning, which can be used to design next-generation algorithms and implementations for competitive, heterogeneous and dynamic cognitive radio systems [45]. The presented framework can serve as a guideline for designing spectrum access solutions that are concerned with the tensions and relationships among autonomous adaptation by secondary (unlicensed) users, the explicit and implicit competition among these users, as well as the interaction of these users with spectrum moderators having their own goals (e.g. making money, imposing fairness rules, ensuring compliance to FCC [1] etc.).

The proposed knowledge-driven framework can be used to design efficient solutions for the usage of the spectrum under a broad set of operating scenarios. These scenarios include “fresh” spectrum, where all radios are cognitive, interactions of cognitive radios with licensed (non-adaptive, high-priority) users, and interactions of cognitive radios with legacy radios in the ISM bands. This framework provides incentives for the secondary users to deploy advanced transmission strategies, to effectively gather information about the environment and learn from on it and, finally, to efficiently share the network resources.

We would like to note though that a large body of research and development work will still need to take place before such a knowledge-driven framework can be deployed. For instance, enhanced learning solutions that make optimal tradeoffs between the resulting utility and implementation costs need to be developed. Moreover, the various solutions for both dynamic spectrum access and learning will need to be tested in heterogeneous and highly dynamic cognitive radio systems, where a variety of SUs are competing for resources. Also, spectrum owners and wireless users will need to decide whether to adopt centralized or distributed solutions for managing the resources, whether they would like to make money, what type of fairness rules they would like to enforce etc.

Finally, we believe that such cognitive radio networking solutions, which are based on stochastic interactions among users rather than the fixed, predetermined solutions and regulations used in the current networks, will ultimately lead to a new generation of cyber-infrastructure, and also next-generation applications, services and intelligent devices. Such solutions are especially necessary to ensure the proliferation of delay-sensitive, high-bandwidth multimedia applications and services, because these are most impacted by the inefficient spectrum use.

## REFERENCES

- [1] Federal Communications Commission, "Spectrum Policy Task Force," *Rep. ET Docket* No. 02-135, Nov. 2002.
- [2] S. Haykin, "Cognitive radio: Brain-empowered wireless communications," *IEEE J. Sel. Areas Commun.*, vol. 23, no. 2, Feb. 2005.
- [3] F. A. Ian, W.Y. Lee, M.C. Vuran, and S. Mohanty, "~~Next~~ generation/dynamic spectrum access/cognitive radio wireless network: a survey," *Computer Networks*, vol. 50, no. 13, Sept. 2006.
- [4] A. Sahai, D. Cabric, N. Hoven, R. Tandra, S.M. Mishra, and R. Brodersen, "Spectrum sensing: fundamental limits and practical challenges," *Tutorial presented at the 2005 DySPAN Conference*.
- [5] A. Ghasemi and E.S. Sousa, "Collaborative spectrum sensing for opportunities access in fading environments," *Proceedings of Dyspan 2005*, pp. 131-136, Nov. 2005.
- [6] C. Doerr, M. Neureld, J. Fifield, et al. "MultiMAC - an adaptive MAC framework for dynamic radio networking," *Proceedings of Dyspan 2005*, pp. 548-555, Nov. 2005.
- [7] C. Kloeck, H. Jaekel, and F. Jondral, "Auction Sequence as a new resource allocation mechanism," *Proceedings of VTC'05*, Dallas, Sept. 2005.
- [8] F. Fu, M. van der Schaar, "Learning for Dynamic Bidding in Cognitive Radio Resources", UCLA Technical Report March 2007 ([available at http://medianetlab.ee.ucla.edu/papers/ffu2007d.pdf](http://medianetlab.ee.ucla.edu/papers/ffu2007d.pdf)).
- [9] J. Huang, R. Berry and M. L. Honig, "Auction-based Spectrum Sharing", *ACM Mobile Networks and Applications Journal (MONET)*, vol. 11, no. 3, pp. 405-418, June 2006.
- [10] L. Berlemann, S. Mangold, G.R. Hiertz and B.H. Walke, "Policy defined spectrum sharing and medium access for cognitive radios", *Journal of Communications, Academy Publishers*, Vol. 1, Issue 1, April 2006.
- [11] C. T. Chou, S. Shankar N, H. Kim and K. Shin, "What and how much to gain by spectrum agility?" *IEEE J. Sel. Areas Commun.*, vol. 25, no. 3, pp. 576-588, Apr. 2007.
- [12] S. Shankar, C.T. Chou, K. Challapali, and S. Mangold, "Spectrum agile radio: capacity and QoS implications of dynamic spectrum assignment," *Global Telecommunications Conference*, Nov. 2005.
- [13] M. van der Schaar, Y. Andreopoulos, Z. Hu, "Optimized scalable video streaming over IEEE 802.11 a/e HCCA wireless networks under delay constraints," *IEEE Trans. Mobile Comput.*, vol. 5, no. 6, pp. 755-768, June 2006.
- [14] "IEEE 802.11e/D5.0, wireless medium access control (MAC) and physical layer (PHY) specifications: Medium access control (MAC) enhancements for Quality of Service (QoS), draft supplement," June 2003.
- [15] F. Fu, T. M. Stoenescu, and M. van der Schaar, "A Pricing Mechanism for Resource Allocation in Wireless Multimedia Applications," *IEEE Journal of Sel. Topics in Signal Process.*, vol. 1, no. 2, pp. 264-279, Aug. 2007
- [16] R. W. Lucky, "Tragedy of the commons," *IEEE Spectrum*, vol. 43, No. 1, pp. 88, Jan 2006.
- [17] R. G. Gallager, "Discrete stochastic processes," Kluwer Academic Publishers, 1996.
- [18] L. S. Shapley, "Stochastic games," *Proceedings of the National Academy of Sciences of the United States of America*, vol. 39, 1095-1100, 1953.
- [19] C. Watkins, and P. Dayan, "Q-learning," *Technical Note, Machine Learning*, vol. 8, 279-292, 1992.
- [20] P. Klemperer, "Auction theory: A guide to the literature," *J. Economics Surveys*, vol. 13, no. 3, pp. 227-286, Jul. 1999.
- [21] J. Sun, E. Modiano, and L. Zheng, "Wireless channel allocation using an auction algorithm," *IEEE J. Sel. Areas Commun.*, vol. 24, no. 5, May 2006.
- [22] F. Fu, and M. van der Schaar, "Non-collaborative resource management for wireless multimedia applications using mechanism design," *IEEE Transaction on Multimedia*, vol. 9, no. 4, pp. 851-868, Jun. 2007.
- [23] D. Fudenberg, and D. K. Levine, "The theory of learning in games," Cambridge, MA: MIT Press, 1999.
- [24] Q. Zhao, L. Tong, A. Swami, and Y. Chen, "~~Decentralized Cognitive MAC for Opportunistic Spectrum Access in Ad Hoc Networks: A POMDP Framework,~~" *IEEE J. Sel. Areas Commun.*, ~~vol. 25, no. 3, April 2007.~~
- [25] M. Jackson, "Mechanism theory," *In the Encyclopedia of Life Support Systems*, 2003.
- [26] Q. Zhang, and S.A. Kassam, "Finite-state Markov model for Rayleigh fading channels," *IEEE Transaction on Communications*, vol. 47, no. 11, Nov. 1999.
- [27] S. Lal, and E.S. Sousa, "Distributed resource allocation for DS-CDMA-based multimedia ad hoc wireless LANs," *IEEE J. Sel. Areas Commun.*, Vol. 17, No. 5, 947 – 967, May 1999.

Deleted: NeXt

Formatted: Font: Italic

Deleted: can be found

Formatted: Font: Italic

Deleted: .

Deleted: ...

- [28] G. Berlanda-Scorza, C. Sacchi, F. Granelli, and F. De Natale, "A QoS-oriented medium access control strategy for variable-bit-rate MC-CDMA transmission in wireless LAN environments," *Proceedings of Globecom*, Vol. 1, 475 – 479, Dec. 2003.
- [29] J. Hu and P. Wellman, "Multiagent reinforcement learning: theoretical framework and an algorithm," *Proceedings of the Fifteenth International Conference on Machine Learning*, pp. 242-250, 1998.
- [30] Y. Su and M. van der Schaar, "A New Perspective on Multi-user Power Control Games in Interference Channels", UCLA Technical Report September 2007 ([available at http://medianetlab.ee.ucla.edu/papers/ysu2007c.pdf](http://medianetlab.ee.ucla.edu/papers/ysu2007c.pdf)).
- [31] Y. Su and M. van der Schaar, "Multi-user Multimedia Resource Allocation over Multi-carrier Wireless Networks," *IEEE Trans. on Signal Process.*, to appear.
- [32] A. Larcher, H. Sun, M. van der Schaar, and Z. Ding, "Decentralized Transmission Strategy for Delay-Sensitive Applications over Spectrum Agile Network," in *Proc. 13th Int. Packet Video Workshop (PV 2004)*, 2004.
- [33] A.B. MacKenzie, L.A. DaSilva, "Game Theory for Wireless Engineers", Synthesis Lectures on Communications, Morgan and Claypool Publishers, 2006.
- [34] H. Ji and C. -Y. Huang, "Non-cooperative uplink power control in cellular radio systems," *Wireless Networks*, vol. 4, no. 3, pp. 233–240, 1998.
- [35] D. Goodman and N. Mandayam, "Power control for wireless data," *IEEE Pers. Comm. Magazine*, vol. 7, no. 2, pp. 48–54, April 2000.
- [36] M.Xiao, N. Schroff, and E. Chong, "Utility based power control in cellular radio systems," in *Proceedings of IEEE INFOCOM*, Anchorage, Alaska, 2001.
- [37] A. B. MacKenzie and S. B. Wicker, "Game theory in communications: Motivation, explanation, and application to power control," in *Proceedings of IEEE GLOBECOM*, pp. 821–826, 2001.
- [38] D. Vengerov, N. Bambos, H. Berenji, "A fuzzy reinforcement learning approach to power control in wireless transmitters," *IEEE Trans. on Systems, Man, and Cybernetics, Part B*, vol. 35, no. 4, pp. 768–778, 2005.
- [39] C. Long, Q. Zhang, B. Li, H. Yang, and X. Guan, "Non-cooperative power control for wireless ad hoc networks with repeated games," *IEEE J. Sel. Areas Commun.*, vol. 25, no. 6, pp. 1101–1112, 2007.
- [40] M. Maskery and V. Krishnamurthy, "Decentralized activation in a zigbee-enabled unattended ground sensor network: A correlated equilibrium game theoretic analysis," in *Proc. of IEEE ICC '07*, 2007, pp. 3915–3920.
- [41] Z. Han, C. Pandana, and K. Liu, "Distributive opportunistic spectrum access for cognitive radio using correlated equilibrium and no-regret learning," *Wireless Communications and Networking Conference (WCNC 2007)*, pp. 11–15, 2007.
- [42] V. Srivastava, J. Neel, A. B. MacKenzie, R. Menon, L. A. DaSilva, J. E. Hicks, J. H. Reed, and R. P. Gilles, "Using game theory to analyze wireless ad hoc networks," *IEEE Communication Surveys and Tutorials*, 2005.
- [43] J. Neel, R. M. Buehrer, J.H. Reed, and R. P. Gilles, "Game theoretic analysis of a network of cognitive radios," in *Proceedings of the 45th Midwest Symposium on Circuits and Systems*, vol. 3, pp. 409–412, August 2002.
- [44] W. Yu, G. Ginis, and J. Cioffi, "Distributed multiuser power control for digital subscriber lines," *IEEE J. Sel. Areas Commun.* vol. 20, no. 5, pp. 1105–1115, June 2002.
- [45] S. Haykin, "Cognitive Dynamic Systems," *Proc. of IEEE ICASSP 2007*, vol. 4, pp. 1369–1372, April 2007.
- [46] Y. Chevaleyre, P. E. Dunne, U. Endriss, J. Lang, M. Lemaître, N. Maudet, J. Padget, S. Phelps, J. A. Rodríguez-Aguilar, and P. Sousa, "Issues in Multiagent Resource Allocation," *Informatica*, 30:3–31, 2006.
- [47] H. P. Young, "Interactive learning and its Limits," Oxford University Press, NY 2004.
- [48] [http://www.infoworld.com/article/06/04/06/77219\\_HNspectrumfrenzy\\_1.html](http://www.infoworld.com/article/06/04/06/77219_HNspectrumfrenzy_1.html)
- [49] F. Meshkati, H.V. Poor, S.C. Schwartz, "Energy-Efficient Resource Allocation in Wireless Networks", *IEEE Signal Proc. Magazine*, Vol. 24, No. 3, pp. 58–68, May 2007.
- [50] H. P. Shiang and M. van der Schaar, "Queueing-Based Dynamic Channel Selection for Heterogeneous Multimedia Applications over Cognitive Radio Networks," *IEEE Trans. Multimedia*, to appear.
- [51] Y.B. Reddy, "Detecting Primary Signals for Efficient Utilization of Spectrum Using Q-Learning," *Proceedings of Fifth International Conference on Information Technology: New Generations (img 2008)*, pp. 360–365, 2008.

Formatted: Font: Italic

Deleted: p

Deleted: 2004

Formatted: Font: Not Italic

Deleted: 2001,

Deleted: IEEE Trans. on,

Deleted: Sel. Areas in Communications

Deleted: IEEE Journal on,

Deleted: in

Deleted: , 2007

Deleted: .

Deleted: IEEE, 2007

Deleted: August 2002,

Deleted: .

Formatted: Font: Not Italic

Formatted: Font: 10 pt

Formatted: Font: Italic

Formatted: Font: (Default) Times New Roman

Formatted: Bullets and Numbering

Formatted: Font: (Default) Times New Roman

Formatted: Font: (Default) Times New Roman

Formatted: Font: 10 pt

Formatted: Font: Italic

Formatted: Font: 10 pt, Italic

Formatted: Font: 10 pt

Formatted: Font: 10 pt