

# Learning to Compete for Resources in Wireless Stochastic Games

Fangwen Fu and Mihaela van der Schaar, *Senior Member, IEEE*

**Abstract**—In this paper, we model the various users in a wireless network (e.g., cognitive radio network) as a collection of selfish autonomous agents that strategically interact to acquire dynamically available spectrum opportunities. Our main focus is on developing solutions for wireless users to successfully compete with each other for the limited and time-varying spectrum opportunities, given experienced dynamics in the wireless network. To analyze the interactions among users given the environment disturbance, we propose a stochastic game framework for modeling how the competition among users for spectrum opportunities evolves over time. At each stage of the stochastic game, a central spectrum moderator (CSM) auctions the available resources, and the users strategically bid for the required resources. The joint bid actions affect the resource allocation and, hence, the rewards and future strategies of all users. Based on the observed resource allocations and corresponding rewards, we propose a best-response learning algorithm that can be deployed by wireless users to improve their bidding policy at each stage. The simulation results show that by deploying the proposed best-response learning algorithm, the wireless users can significantly improve their own bidding strategies and, hence, their performance in terms of both the application quality and the incurred cost for the used resources.

**Index Terms**—Delay-sensitive transmission, interactive learning, multiuser resource management, reinforcement learning, stochastic games, wireless networks.

## I. INTRODUCTION

**D**YNAMIC resource management in heterogeneous wireless networks is a challenging problem [3]. The wireless stations and radio systems that must coexist in such a network differ in their individual utility functions, transmission actions, resource demands, and capabilities. Thus, various levels of strategic<sup>1</sup> interaction and adaptation are necessary to cope with the widely varying dynamics. In this paper, we focus on synthesizing new, dynamic, and informationally decentralized resource-management mechanisms to achieve high utility in competitive and heterogeneous wireless networks (including cognitive radio networks [1]–[3]). Specifically, our focus is on designing associated communication algorithms that enable

self-interested autonomous wireless stations to strategically compete for the available spectrum resources in either ISM bands [1] or bands shared with licensed users, according to *a priori* mandated or negotiated rules.

This paper is primarily concerned with the tensions and relationships among autonomous adaptation by secondary (unlicensed) users (SUs), the competition among these users, the interaction of these users with spectrum moderators having their own goals, e.g., making money, imposing fairness rules, ensuring compliance with the Federal Communications Commission (FCC) [1], and local regulations with respect to primary (licensed) users (PUs), etc. Unlike previous works on resource management [6], [21], [26], our main focus is on discussing how users can adapt, predict, learn, and determine how they compete for the time-varying resources, as well as how they select the associated transmission strategies, given the experienced “dynamics.”

In wireless networks, these dynamics can be categorized into two types: One is the *disturbance due to the “environment,”* and the other is the *impact caused by competing users.* The disturbance due to the environment results from variations (uncertainties) of the wireless channels or source (e.g., multimedia) characteristics. For example, the stochastic behavior of the PUs, the time-varying channel conditions experienced by the SUs, and the time-varying source traffic that needs to be transmitted by the SUs can be considered as environmental disturbances. These types of dynamics are generally modeled as stationary processes. For instance, the use of each channel by the PUs can be modeled as a two-state Markov chain with ON-state (the channel is used by PUs) and OFF-state (the channel is available for the SUs) [7]. The channel conditions can be modeled using a finite-state Markov model [24]. The packet arrival of the source traffic can be modeled as a Poisson process<sup>2</sup> [11].

Conventionally, wireless stations have only considered these environment disturbances when adapting their cross-layer strategies [12] for delay-sensitive transmission. The other type of dynamics—the impact from competing users, which is due to the noncollaborative, autonomous, and strategic SUs in the network transmitting their traffic—is less well studied to wireless communication networks.

The goal of this paper is to provide solutions and associated metrics that can be used by an autonomous SU to analyze and predict the outcome of various dynamic interactions among competing SUs in dynamic multiuser communication

Manuscript received August 28, 2007; revised April 17, 2008 and July 1, 2008. This work was supported by the National Science Foundation under CAREER Award CCF-0541867. The review of this paper was coordinated by Prof. O. B. Akan.

The authors are with the Department of Electrical Engineering, University of California at Los Angeles, Los Angeles, CA 90095 USA (e-mail: ffwu@ee.ucla.edu; mihaela@ee.ucla.edu).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TVT.2008.2002917

<sup>1</sup>By strategic users, we mean users that are not price takers and do not have an *a priori* consensus on resource allocation.

<sup>2</sup>Other packet arrival models can also be used in our proposed framework.

87 systems and, based on this forecast, adapt and optimize its  
 88 transmission strategy. In our considered wireless networks,  
 89 the SUs are modeled as rational and strategic. We model the  
 90 spectrum management as a stochastic game [22] in which the  
 91 SUs simultaneously and repeatedly make their own resource  
 92 bids. The competition for dynamic resources is assisted by a  
 93 central coordinator (similar to that in existing wireless LAN  
 94 (WLAN) standards such as 802.11e HCF [13]). We refer to this  
 95 coordinator as the central spectrum moderator (CSM). The role  
 96 of the CSM is to allocate resources to the SUs based on the  
 97 predetermined utility maximization rule.<sup>3</sup>

AQ1

98 In this paper, to explicitly consider the strategic behavior of  
 99 autonomous SUs and the informationally decentralized nature  
 100 of the competition for wireless resources, we assume that the  
 101 CSM deploys an auction mechanism for dynamically allocat-  
 102 ing resources. Auction theory has extensively been studied  
 103 in economics [19], and it has also been recently applied to  
 104 network resource allocation [4]–[6]. Note that the role of the  
 105 CSM<sup>4</sup> in our resource management game for our considered  
 106 wireless networks will be kept to a minimum. Unlike alternative  
 107 existing solutions [21], the CSM will not require knowledge  
 108 of the private information of the users and will not perform  
 109 complex computations for deciding the resource allocation. Its  
 110 only role will be the implementation of the spectrum etiquette  
 111 rules as in [8] and ensuring that the available spectrum holes  
 112 are auctioned among users. To capture the network dynamics,  
 113 we allow the CSM to *repeatedly* auction the available spectrum  
 114 opportunities based on the PUs' behaviors. Meanwhile, each  
 115 SU is allowed to strategically adapt its bidding strategy based  
 116 on information about the available spectrum opportunities, its  
 117 source and channel characteristics, and the impact of the other  
 118 SU bidding actions.

119 Using this stochastic wireless allocation framework, we de-  
 120 velop a learning methodology for SUs to improve their policies  
 121 for playing the auction game, i.e., the policies for generating  
 122 the bids to compete for available resources. Specifically, during  
 123 repeated multiuser interaction, the SUs can observe partial his-  
 124 toric information of the outcome of the auction game, through  
 125 which the SUs can estimate the impact on their future rewards  
 126 and then adopt their best response to effectively compete for  
 127 channel opportunities. The estimation of the impact on the  
 128 expected future reward can be performed using different types  
 129 of interactive learning [18]. In this paper, we focus on reinforce-  
 130 ment learning [17], [27] because this allows the SUs to improve  
 131 their bidding strategy based only on the knowledge of their own  
 132 past received payoffs without knowing the bids or payoffs of  
 133 the other SUs. Our proposed best-response learning algorithm  
 134 is inspired from the Q-learning for the single agent interact-  
 135 ing with the environment. Unlike Q-learning, the proposed  
 136 best-response learning explicitly considers the interactions and  
 137 coupling among SUs in the wireless network. By deploying  
 138 the best-response learning algorithm, the SUs can strategically

predict the impact of current actions on future performance and  
 then optimally make their resource bids.

This paper is organized as follows. In Section II, we intro-  
 duce a stochastic game formulation for multiuser interaction  
 in wireless networks. In Section III, we show how a one-  
 stage auction mechanism can be used to divide the spectrum  
 allocation among strategic SUs. In Section IV, we present  
 the state definition, state transition model, and stage reward  
 function for the SUs in the stochastic game. In Section V,  
 we discuss the bidding strategies of the SUs for playing the  
 stochastic game. In Section VI, we propose a best-response  
 learning approach for the SUs to predict their future rewards  
 based on the observed historic information. In Section VII,  
 we present the simulation results, followed by conclusions and  
 future research in Section VIII.

## II. STOCHASTIC GAME FORMULATION FOR DYNAMIC MULTIUSER INTERACTION

We consider a spectrum consisting of  $N$  channels, each  
 indexed by  $j \in \{1, \dots, N\}$ . The  $N$  wireless channels are orig-  
 inally licensed to a primary network (PN) whose users (i.e.,  
 PUs) exclusively access the channels. In the secondary network  
 (SN), the  $M$  ( $M \geq N$ ) autonomous SUs, each indexed by  
 $i \in \{1, \dots, M\}$  and transmitting delay-sensitive data, compete  
 for the spectrum opportunities released by the PUs in these  
 $N$  channels. Although the available transmission opportunities  
 (TxOps) for SUs depend on the access patterns of PUs and the  
 detection systems [2], we do not discuss the detection methods  
 in this paper but rather rely on the existing literature for this  
 purpose [3]. Instead, we assume that the available TxOps in  
 each channel change over time due to the PUs joining or leaving  
 the network and can be modeled as a two-state Markov chain,  
 as in [7] and [10]. Our goal is to develop a general framework  
 for multiuser interaction in the SN, where users can compete  
 for dynamically available TxOps. Moreover, we also aim to  
 provide solutions for SUs to improve their strategies for playing  
 the repeated resource-management game by considering their  
 past interactions with other SUs.

The communications of the PUs are assumed to follow a  
 synchronous slot structure. The time slot has length of  $\Delta T$   
 seconds. We assume that during each time slot, each channel  
 is either exclusively occupied by PUs or that there is no PU  
 accessing the channel [7], [10]. Hence, during each time slot,  
 the channel is in one of the following two states: ON-state  
 (this channel is currently used by the PUs) or OFF-state (this  
 channel is not used by the PUs, and hence, the SUs can use this  
 channel). Note that if this is an unlicensed band, the channel  
 will always be in the off mode and can be utilized by the  
 SUs at all times. The TxOp of channel  $j$  at time slot  $t \in \mathbb{N}$   
 is denoted by  $y_j^t \in \{0, 1\}$ , where  $y_j^t$  is 0 if the channel is  
 in the ON-state and 1 if it is in the OFF-state. In this paper,  
 the TxOp  $y_j^t$  of channel  $j$  is modeled by a two-state Markov  
 chain with transition probability  $p_j^{FN} = p(y_j^{t+1} = 0 | y_j^t = 1)$   
 and  $p_j^{NF} = p(y_j^{t+1} = 1 | y_j^t = 0)$ . The TxOp profile of the  
 $N$  channels is represented by  $\mathbf{y}^t = [y_1^t, \dots, y_N^t]$ .

As in [13], we assume that a polling-based medium-access  
 protocol is deployed in the SN, which is arbitrated by a CSM.

<sup>3</sup>Other fairness rules can also be deployed in the CSM such as air-time fairness, utility-based fairness, etc. [12].

<sup>4</sup>It should be noted that this approach can also allow for multiple CSMs to manage the spectrum by fairly dividing their responsibilities, e.g., based on their geolocation or frequency band in which they are operating, or by competing against each other for the number of SUs that will associate with them.

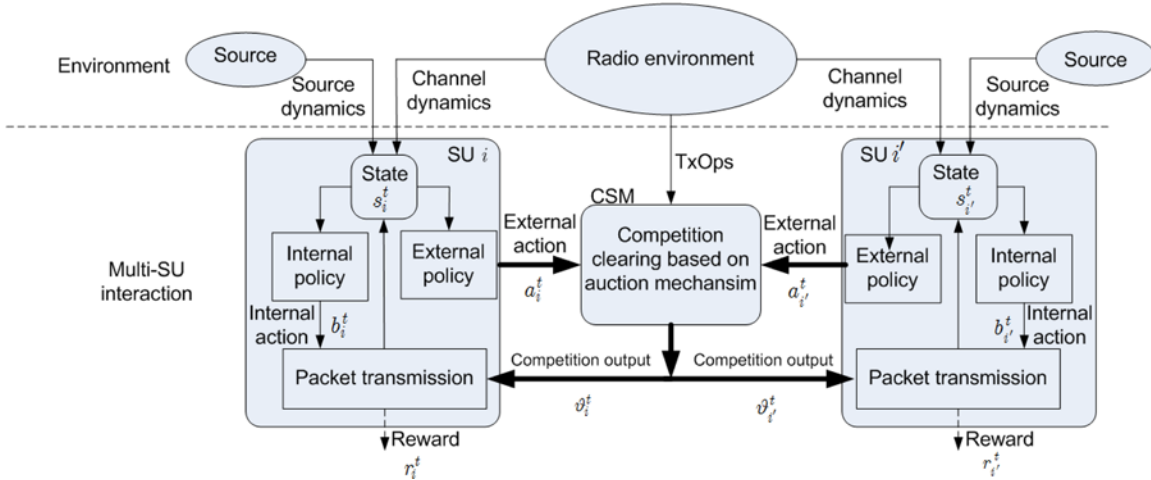


Fig. 1. Conceptual overview of the multi-SU interaction in the SN.

195 The polling policy is only changed at the start of every time  
 196 slot. For simplicity, we assume that each SU can access a  
 197 single channel, and that each channel can be accessed by  
 198 a single SU within the time slot. The SUs can switch the  
 199 channels only when crossing time slots. Note that this simple  
 200 medium-access model used for illustration in this paper can  
 201 easily be extended to more sophisticated models [10], where  
 202 each SU can simultaneously access multiple channels or the  
 203 channels are being shared by multiple SUs, etc. When using  
 204 this time-division channel access, we assume that the wireless  
 205 users deploy constant transmission power and experience no  
 206 interference. Furthermore, we assume that the wireless users  
 207 move slowly, and thus, their experienced channel conditions  
 208 slowly change.

209 During each time slot, an SU needs to first determine how to  
 210 compete with the other SUs for the time-varying TxOps. This  
 211 represents its external actions, since they determine the inter-  
 212 action between this SU and the other SUs, and the amount of  
 213 resources allocated to that SU. The external actions at time slot  $t$   
 214 are denoted by  $a_i^t \in A_i$ , where  $A_i$  is the set of possible external  
 215 actions available to SU  $i$ . Based on the allocated resources,  
 216 the SU determines how to transmit its traffic (application layer  
 217 data) by selecting the various strategies at different layers of  
 218 the OSI stack (e.g., through cross-layer adaptation [12]). These  
 219 actions are referred to as internal actions, since they only  
 220 determine the SU's utility at the current time. The internal  
 221 actions at time slot  $t$  are denoted by  $b_i^t \in B_i$ , where  $B_i$  is the set  
 222 of possible internal actions available to SU  $i$ . In this paper, we  
 223 propose an auction mechanism deployed in the CSM. Hence,  
 224 the external action  $a_i^t$  of SU  $i$  is the bid it submits to CSM. The environ-  
 225 ment experienced by an SU  $i$  can be characterized by its current  
 226 "state"  $s_i^t \in S_i$ , which will be discussed in Section IV. At each  
 227 time slot  $t$ , SU  $i$  generates the external action  $a_i^t$  to compete  
 228 for the TxOps  $\mathbf{y}^t$ . The competition result is  $\vartheta_i^t$ , based on which  
 229 SU  $i$  performs its internal action  $b_i^t$  and obtains the reward  $r_i^t$  at  
 230 this time slot. After packet transmission, SU  $i$  transits to the  
 231 next state  $s_i^{t+1} \in S_i$ . The conceptual overview of the multi-  
 232 SU interactions in the repeated auctions is illustrated in Fig. 1.

234 The repeated competition among the SUs can be modeled as  
 235 a stochastic game [16], [22]. The time slot corresponds to the  
 236 term "stage," which is commonly used in stochastic games. In  
 237 the remainder of this paper, we interchangeably use the terms  
 238 "time slot" and "stage."

239 We define the stochastic game for SN resource allocation as  
 240  $\langle \{S_i, A_i, B_i, O_i, q_i, r_i\}_{i=1}^M, \mathcal{Y} \rangle$ , where each SU  $i$  is associated  
 241 with a tuple  $\langle S_i, A_i, B_i, O_i, q_i, r_i \rangle$ . Specifically, we have the  
 242 following.

- 1)  $\mathcal{Y}$  is a finite set of possible TxOps available for SUs. 243  
 In this paper,  $\mathcal{Y} = \{0, 1\}^N$ , and  $\mathbf{y}^t \in \mathcal{Y}$  is the avail- 244  
 able TxOps at stage  $t$ , which is common information 245  
 for SUs. 246
- 2)  $S_i$  is a finite local state space of SU  $i$ . We let  $S := 247$   
 $\prod_{k=1}^N S_k$  be the global state space of all SUs and 248  
 $S_{-i} := \prod_{k \neq i} S_k$  be the global state space of SUs other 249  
 than  $i$ . At stage  $t$ , the global state is denoted by  $\mathbf{s}^t = 250$   
 $(s_1^t, \dots, s_M^t) = (s_i^t, \mathbf{s}_{-i}^t)$ , where  $-i$  represents all the 251  
 SUs other than  $i$ . 252
- 3)  $A_i$  is a finite set of external actions performed by SU  $i$  253  
 to compete for the available TxOps. The external action 254  
 vector at stage  $t$  for all SUs will be  $\mathbf{a}^t = (a_1^t, \dots, a_M^t)$ . 255
- 4)  $B_i$  is a finite set of internal actions performed by SU  $i$  to 256  
 determine the packet transmission. 257
- 5)  $O_i$  is a finite set of possible output from multi-SU com- 258  
 petition. In this paper, the output  $\vartheta_i^t \in O_i$  is the auction 259  
 result computed by the CSM for SU  $i$  at stage  $t$ . We will 260  
 give the specific form of the output in Section III. 261
- 6)  $q_i$  is the state transition probability for SU  $i$ . Thus, 262  
 $q_i(s_i^{t+1}, \mathbf{y}^{t+1} | s_i^t, \mathbf{y}^t, \vartheta_i^t, b_i^t)$  is the probability that the state 263  
 of SU  $i$  transits from  $s_i^t$  to  $s_i^{t+1}$  and TxOp transits from 264  
 $\mathbf{y}^t$  to  $\mathbf{y}^{t+1}$  if the competition output is  $\vartheta_i^t$  and the internal 265  
 action is  $b_i^t$ . The reason that the transition probability 266  
 includes the common TxOp  $\mathbf{y}^t$  is because the channel 267  
 condition transition of SU  $i$  depends on the available 268  
 TxOp. 269
- 7)  $r_i$  is the stage reward (immediate reward) received by SU 270  
 $i$ , where  $r_i : (S_i, O_i, B_i) \mapsto \mathbb{R}$ . It should be noted that 271

272 the reward function  $r_i$  depends on the competition output  
 273 and, hence, indirectly depends on the other SUs' external  
 274 actions.

275 To design a stochastic game for the SN with strategic SUs,  
 276 we have to consider the following: 1) What auction mech-  
 277 anism can be deployed to resolve the competition among  
 278 SUs; 2) how the dynamic environment experienced by each  
 279 SU can be modeled; and 3) how the SUs can forecast the  
 280 impact of their bids made at the current time on their future  
 281 performance?

### 282 III. AUCTION MECHANISM—ONE STAGE 283 RESOURCE ALLOCATION

284 In this paper, we assume that the CSM is aware of the  
 285 TxOp  $\mathbf{y}^t$  and allocates (through polling the SUs) those channels  
 286 with  $y_j^t = 1$  to the SUs. To efficiently allocate the available  
 287 resources (opportunities), the CSM needs to collect information  
 288 about the SUs [21]. However, as mentioned in Section I, in a  
 289 wireless network, the information is decentralized, and thus,  
 290 the information exchange between the SUs and the CSM needs  
 291 to be kept limited due to the incurred communication cost.  
 292 On the other hand, the SUs competing with each other are  
 293 selfish and strategic, and hence, the information they hold is  
 294 private, and they may not desire to reveal this information to  
 295 the CSM. Therefore, one of our key interests in this paper is  
 296 to determine what information should be exchanged between  
 297 the SUs and the CSM and how this information should be  
 298 exchanged. In the following, we present an auction mechanism  
 299 for dynamically coordinating the interactions among SUs and  
 300 discuss the computational complexity in the CSM and the  
 301 communication cost between SUs and CSM.

302 First, the CSM announces the auction by broadcasting the  
 303 TxOp  $\mathbf{y}^t$ . The SUs receive the announcement and determine the  
 304 external action (i.e., the bid vector)  $\mathbf{a}_i^t = [a_{i1}^t, \dots, a_{iN}^t] \in \mathbb{R}^N$   
 305 based on the announced information and their own private  
 306 information about the environment they experience, which is  
 307 discussed in detail in Section IV. Subsequently, each SU sub-  
 308 mits the bid vector to the CSM. After receiving the bid vectors  
 309 from the SUs, the CSM computes the channel allocation  $\mathbf{z}_i^t =$   
 310  $[z_{i1}^t, \dots, z_{iN}^t] \in \{0, 1\}^N$  for each SU  $i$  based on the submitted  
 311 bids. To compel the SUs to truthfully declare their bids [23],  
 312 the CSM also computes the payment  $\tau_i^t \in \mathbb{R}_-$  that the SUs have  
 313 to pay for the use of resources during the current stage of the  
 314 game. The negative value of the payment means the absolute  
 315 value that SU  $i$  has to pay the CSM for the used resources.  
 316 Hence, the competition output  $\vartheta_i^t$  in this auction mechanism  
 317 includes the channel allocation  $\mathbf{z}_i^t$  and the payment  $\tau_i^t$ , i.e.,  
 318  $\vartheta_i^t = (\mathbf{z}_i^t, \tau_i^t)$ . The competition output is then transmitted back  
 319 to the SUs. The computation of the channel allocation  $\mathbf{z}_i^t$  and  
 320 payment  $\tau_i^t$  is described as follows.

321 After each SU submits the bid vector, the CSM performs  
 322 two computations, i.e., channel allocation and payment com-  
 323 putation. Note that most existing multiuser wireless resource  
 324 allocation solutions can be modeled as such repeated auctions  
 325 for resources. If the resources are priced or the users may lie  
 326 about their resource needs, taxes associated with the resource

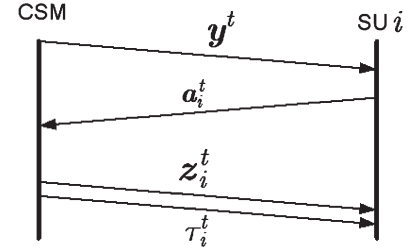


Fig. 2. Information exchange between CSM and SU  $i$ .

usage will need to be imposed [14]. Otherwise, these taxes can  
 be considered to be zero throughout the paper.

We denote the channel allocation matrix  $Z^t = [z_{ij}^t]_{M \times N}$   
 with  $z_{ij}^t$  being 1 if channel  $j$  is assigned to SU  $i$ , and 0  
 otherwise. The feasible set of channel assignments is denoted  
 as  $\mathcal{Z}^t = \{Z^t | \sum_{i=1}^M z_{ij}^t = y_j^t, \forall j, \sum_{j=1}^N z_{ij}^t \leq 1, \forall i, z_{ij}^t \in$   
 $\{0, 1\}\}$ . The channel allocation matrix without the pres-  
 ence of SU  $i$  is denoted  $Z_{-i}^t = [z_{kj}^t]_{(M-1) \times N}$ , and the  
 corresponding feasible set is  $\mathcal{Z}_{-i}^t = \{Z_{-i}^t | \sum_{k=1, k \neq i}^M z_{kj}^t =$   
 $y_j^t, \forall j, \sum_{j=1}^N z_{kj}^t \leq 1, \forall k \neq i, z_{kj}^t \in \{0, 1\}\}$ , where  $-i =$   
 $\{1, \dots, i-1, i+1, \dots, M\}$ . During the first phase, the CSM  
 allocates the channels to SUs based on its adopted fairness rule,  
 e.g., maximizing the total “social welfare,”<sup>5</sup> as

$$Z^{t, \text{opt}} = \arg \max_{Z^t \in \mathcal{Z}^t} \sum_{i=1}^M \sum_{j=1}^N z_{ij}^t a_{ij}^t. \quad (1)$$

If the resources are priced, we will consider in this paper,  
 for illustration, a second-price auction mechanism [19], [23] for  
 determining the tax that needs to be paid by SU  $i$  based on the  
 above optimal channel assignment  $Z^{t, \text{opt}} = [z_{ij}^{t, \text{opt}}]_{M \times N}$ . This  
 tax is equal to

$$\tau_i^t = \sum_{k=1, k \neq i}^M \sum_{j=1}^N z_{kj}^{t, \text{opt}} a_{kj}^t - \max_{Z_{-i}^t \in \mathcal{Z}_{-i}^t} \sum_{k=1, k \neq i}^M \sum_{j=1}^N z_{kj}^t a_{kj}^t. \quad (2)$$

Note that when  $N = 1$ , the generalized auction mechanism  
 presented above becomes the well-known second-price auction  
 [19]. Although the optimization problems in (1) and (2) are  
 discrete optimizations, they can efficiently be solved using  
 linear programming. As argued in [20], the linear optimization  
 problem can be solved in polynomial time, and hence, the CSM  
 only requires limited computational complexity.

The information exchange between the CSM and the SUs  
 is illustrated in Fig. 2. From Fig. 2, we note that, at each  
 stage, the CSM first broadcasts the available TxOps to all the  
 SUs for the auction, and then each SU submits its own bid  
 vector over all the available TxOps. After receiving the bids, the  
 CSM computes the auction results and sends back to the users  
 the channel allocations and the corresponding payments. The  
 signaling required for the auction is most often implemented  
 at the application layer. In the worst case, the amount of

<sup>5</sup>Note that other fairness solutions than maximizing the social welfare could  
 be adopted, and this will not influence our proposed solution.

361 data communicated between the CSM to the SUs is equal to  
 362  $(M + 1)N + nN$  bits, where  $n$  is the amount of bits repre-  
 363 senting the payment for each SU. The amount of data commu-  
 364 nicated by each SU to the CSM is  $n'N$  bits, where  $n'$  is the  
 365 amount of bits representing the bid submitted to the CSM on  
 366 each channel.

367 Compared with traditional one-stage resource allocation  
 368 methods, our proposed auction mechanism has the following  
 369 advantages.

- 370 1) Unlike traditional centralized resource allocation meth-  
 371 ods [30], our proposed auction mechanism is not required  
 372 to know the SUs' utility functions or preferences, which  
 373 is often the private information of the users and is not  
 374 common knowledge. In fact, our auction mechanism only  
 375 requires the SUs to submit their bid vectors for the avail-  
 376 able TxOps. The bid vector computation is performed  
 377 by the SUs, but not the CSM, based on their utili-  
 378 ties, preferences, action sets, experienced environment  
 379 characteristics, etc.
- 380 2) Unlike traditional decentralized resource allocation meth-  
 381 ods [28] where multiple iterations are required before  
 382 convergence, our proposed auction mechanism only re-  
 383 quires the SUs to submit the bid vectors once. Hence,  
 384 our proposed auction mechanism is suitable for online  
 385 resource management. Moreover, we do not assume as in  
 386 [29] that users are price takers and that there is consensus  
 387 about what is a fair distribution of the resources. Instead,  
 388 in the proposed framework, users are strategic and are  
 389 able to determine their own bid vectors for resources  
 390 based on their knowledge, utilities, preferences, etc.

#### 391 IV. USER MODELING IN THE STOCHASTIC 392 GAME FRAMEWORK

##### 393 A. Definition of SU States

394 As discussed in Section I, each SU needs to cope with two  
 395 types of "uncertainties," i.e., disturbances from the environment  
 396 and interactions with other SUs. The environment is charac-  
 397 terized by packet arrivals from the source (i.e., source/traffic  
 398 characterization) connected with the transmitter and the chan-  
 399 nel conditions. In this section, we will illustrate how these  
 400 disturbances can be modeled. However, note that other models  
 401 of the environment existing in the literature can be adopted. The  
 402 use of a specific model will only affect the performance of the  
 403 proposed solution and not the general framework for multiuser  
 404 interaction proposed in this paper.

405 For illustration, we assume that each SU  $i$  maintains a buffer  
 406 with limited size  $X_i$ , which can be interpreted as a time window  
 407 that specifies which packets are considered for transmission at

each time based on their delay deadlines. Expired packets are  
 408 dropped from the buffer. This model has extensively been used  
 409 for delay-sensitive data transmission, e.g., leaky bucket model  
 410 for video transmission [25]. The number of packets in the buffer  
 411 at time slot  $t$  is denoted as  $x_i^t$  ( $0 \leq x_i^t \leq X_i$ ). We assume that  
 412 the packets arrive from the source at the beginning of each time  
 413 slot, i.e.,  $x_i^t$  is only updated at the beginning of a time slot. The  
 414 number of packets arriving into the buffer during one time slot  
 415 is a random variable independent of the time  $t$  and denoted as  
 416  $\chi_i$ .  $\chi_i$  follows the Poisson distribution with the average arrival  
 417 rate  $\bar{\chi}_i$  packets per second [11]. However, note that the Poisson  
 418 process is simply used for illustration purposes, and other traffic  
 419 models (e.g., renewal process, etc.) can also be used in our  
 420 framework. The average number of packets arriving during one  
 421 time slot is equal to  $\bar{\chi}_i \Delta T$  [11]. 422

The condition of channel  $j$  experienced by SU  $i$  is rep-  
 423 resented by the signal-to-noise ratio (SNR) and denoted as  
 424  $\rho_{ij}^t$  (in decibels). When  $y_j^t = 1$ , we assume that the channel  
 425 condition of each channel can be represented by a set of discrete  
 426 SNR values, i.e.,  $\rho_{ij}^t \in \{\sigma_{ij}^1, \dots, \sigma_{ij}^K\}$ . Note that the number of  
 427 discrete SNR values  $K$  can be determined by SU  $i$  by trading  
 428 off the complexity (a larger  $K$  leads to a larger state space) and  
 429 the resulting impact on performance. When  $y_j^t = 0$ , we set  $\rho_{ij}^t$   
 430 equal to  $-\infty$ , which means that the channel is unavailable to  
 431 SUs at that time. As shown in [24], when  $y_j^t = 1$ , the channel  
 432 condition (in terms of SNR) can also be modeled as a finite-state  
 433 Markov chain, where the transition from channel condition  $\sigma_{ij}^l$   
 434 at time  $t$  to channel condition  $\sigma_{ij}^k$  at time  $t + 1$  takes place with  
 435 probability  $p_{ij}^{l \rightarrow k}$ . These transition probabilities can easily be  
 436 estimated by SU  $i$  by repeatedly interacting with the channel.  
 437 We denote by  $p_{ij}^{-\infty \rightarrow k}$  the probability that the channel condi-  
 438 tion is  $\sigma_{ij}^k$  at time  $t + 1$ , knowing that  $y_j^t = 0$  and  $y_j^{t+1} = 1$ .  
 439 The probability that the channel condition transition to  $-\infty$ ,  
 440 knowing that  $y_j^{t+1} = 0$ , is 1 no matter in what condition the  
 441 channel  $j$  is at time  $t$ . Then, the combination  $(y_j^t, \rho_{ij}^t)$  is still a  
 442 Markov chain with state transition probability as in (3), shown  
 443 at the bottom of the page. 444

To model the dynamics experienced by SU  $i$  at time  $t$  in  
 445 the SN, we define a "state"  $s_i^t = (v_i^t, \rho_i^t) \in \mathcal{S}_i$ , where  $\rho_i^t =$   
 446  $(\rho_{i1}^t, \dots, \rho_{iN}^t)$ . The state encapsulates the current buffer state  
 447 as well as the state of each channel.  $\mathcal{S}_i$  is the set of possible  
 448 states.<sup>6</sup> The total number of possible states for SU  $i$  is equal to  
 449  $|\mathcal{S}_i| = (X_i + 1) \times (K + 1)^N$ . We will show later in this paper  
 450 that the state information is sufficient for SU  $i$  to compete for  
 451 resources (make bid vector) at the current time. 452

<sup>6</sup>We assume that the channel state and the transmission buffer independently evolve as time goes by.

$$p(y_j^{t+1}, \rho_{ij}^{t+1} | y_j^t, \rho_{ij}^t) = \begin{cases} (1 - p_j^{FN}) p_{ij}^{l \rightarrow k}, & \text{if } y_j^t = 1, \quad \rho_{ij}^t = \sigma_{ij}^l, \quad y_j^{t+1} = 1, \quad \rho_{ij}^{t+1} = \sigma_{ij}^k \\ p_j^{NF} p_{ij}^{-\infty \rightarrow k}, & \text{if } y_j^t = 0, \quad y_j^{t+1} = 1, \quad \rho_{ij}^{t+1} = \sigma_{ij}^k \\ p_j^{FN}, & \text{if } y_j^t = 1, \quad \rho_{ij}^t = \sigma_{ij}^l, \quad y_j^t = 0 \\ 1 - p_j^{NF} & \text{o. w.} \end{cases} \quad (3)$$

### 453 B. State Transition and Stage Reward

454 We will now discuss the state transition process. Remember  
 455 that the state of SU  $i$  includes the buffer state  $v_i^t$  and the  
 456 channel state  $\rho_i^t$ . In this paper, we assume that the channel  
 457 state transition is independent of the buffer state transition.  
 458 In the above, we describe the transition of the channel state  
 459  $\rho_i^t$  and the TxOp  $\mathbf{y}^t$ . The buffer state transition is determined  
 460 by the number of packets arriving and the channel allocation  
 461  $z_i^t$  as well as the internal action  $b_i^t$  during that time slot.  
 462 The number of packets transmitted at stage  $t$  is denoted by  
 463  $\mathcal{N}_i(s_i^t, z_i^t, b_i^t)$ . Given the channel allocation, SU  $i$  can adapt  
 464 its own internal action to maximize the number of transmitted  
 465 packets, i.e.,

$$n_i(s_i^t, z_i^t) = \max_{b_i^t \in B_i} \mathcal{N}_i(s_i^t, z_i^t, b_i^t). \quad (4)$$

466 The optimization can be performed by a cross-layer adaptation  
 467 algorithm as in [5], [12], and [21]. Since our focus is on the  
 468 multi-SU interaction, we assume that the internal action will  
 469 always be performed to maximize the number of transmitted  
 470 packets. We simply use  $n_i(s_i^t, z_i^t)$  to represent the number  
 471 of transmitted packets and omit the internal actions in the  
 472 following notations.

473 The evolution of the buffer state is captured by  
 474  $v_i^{t+1} = \min\{(v_i^t - n(s_i^t, z_i^t))^+ + \chi_i^t, X_i\}$ . We define  $h = v_i^{t+1} -$   
 475  $(v_i^t - n(s_i^t, z_i^t))^+$ . Based on the packet arrival model, the buffer  
 476 state transition probability is computed as in (5), shown at the  
 477 bottom of the page. The state transition combined with TxOps,  
 478 given the current resource allocation  $z_i^t$ , can be computed as

$$q_i(s_i^{t+1}, \mathbf{y}^{t+1} | s_i^t, \mathbf{y}^t, z_i^t) \\ = \underbrace{p_i^{\text{buf}}(v_i^{t+1} | v_i^t, z_i^t)}_{\text{buffer state transition}} \prod_{j=1}^N \underbrace{p(y_j^{t+1}, \rho_{ij}^{t+1} | y_j^t, \rho_{ij}^t)}_{\text{channel state transition}} \quad (6)$$

479 where the first term represents the buffer state transition, which  
 480 is independent of the second term of the channel state transition.

481 Based on the channel allocation  $z_i^t$ , the SU transmits  
 482 the available packets in the buffer. In the next time slot,  
 483 new packets arrive into the buffer. Newly incoming packets  
 484 may lead to packets already existing in the buffer being  
 485 dropped whenever the buffer is full or their delay dead-  
 486 line has passed. Clearly, the performance of the application  
 487 (e.g., video quality) improves when fewer packets are lost.  
 488 Hence, we can interpret a negative value of the number of  
 489 lost packets as the stage gain, which is denoted by  $g_i^t$ , i.e.,  
 490  $g_i^t(s_i^t, z_i^t) = -((v_i^t - n_i(s_i^t, z_i^t))^+ + \chi_i^t - X_i)^+$ . The reward at  
 491 time  $t$  for SU  $i$  is expressed using the quasi-linear form  
 492  $r_i(s_i^t, \vartheta_i^t) = g_i^t + \tau_i^t$ . Note that the gain  $g_i^t$  and payment  $\tau_i^t$

depend on the states and bids of all the competing SUs in the 493  
 SN. Hence, the reward is also rewritten as  $r_i(s^t, \mathbf{y}^t, \mathbf{a}^t)$ . 494

## V. BIDDING STRATEGY FOR PLAYING THE STOCHASTIC GAME 495

### A. Best-Response Bidding Policy 497

In the SN, we assume that the stochastic game is played 498  
 by all the SUs for an infinite number of stages. This 499  
 assumption is reasonable for applications having a long 500  
 duration, such as video streaming. In our network setting, we 501  
 define a history of the stochastic game up to time  $t$  as  $\mathbf{h}^t =$  502  
 $\{s^0, \mathbf{y}^0, \mathbf{a}^0, z^0, \tau^0, \dots, s^{t-1}, \mathbf{y}^{t-1}, \mathbf{a}^{t-1}, z^{t-1}, \tau^{t-1}, s^t, \mathbf{y}^t\} \in$  503  
 $\mathcal{H}^t$ , which summarizes all previous states, available TxOps, 504  
 and the actions taken by the SUs as well as the outcomes at 505  
 each stage of the auction game, and  $\mathcal{H}^t$  is the set of all possible 506  
 histories up to time  $t$ . However, during the stochastic game, 507  
 each SU  $i$  cannot observe the entire history but rather part of 508  
 the history  $\mathbf{h}^t$ . The observation of SU  $i$  is denoted as  $\mathcal{O}_i^t \in \mathcal{O}_i^t$  509  
 and  $\mathcal{O}_i^t \subset \mathbf{h}^t$ . Note that the current state  $s_i^t$  can always be 510  
 observed, i.e.,  $s_i^t \in \mathcal{O}_i^t$ . In this paper, we focus on the external 511  
 action selection for the SUs. The external action selection 512  
 for SU  $i$  to play the stochastic game is also referred to as a 513  
 bidding policy  $\pi_i^t: \mathcal{O}_i^t \mapsto A_i$  for SU  $i$  at time  $t$  and defined 514  
 as a mapping from the observations up to the time  $t$  into the 515  
 specific action, i.e.,  $\mathbf{a}_i^t = \pi_i^t(\mathcal{O}_i^t)$ . Furthermore, a policy profile 516  
 $\boldsymbol{\pi}_i$  for SU  $i$  aggregates the bidding policies about how to play 517  
 the game over the entire course of the stochastic game, i.e., 518  
 $\boldsymbol{\pi}_i = (\pi_i^0, \dots, \pi_i^t, \dots)$ . The policy profile for all the SUs at 519  
 time slot  $t$  is denoted as  $\boldsymbol{\pi}^t = (\pi_1^t, \dots, \pi_M^t) = (\boldsymbol{\pi}_i^t, \boldsymbol{\pi}_{-i}^t)$ . 520

The policy  $\boldsymbol{\pi}_i$  is said to be Markov if the bidding policy 521  
 $\pi_i^t$  for  $\forall t$  is, given the current state  $s_i^t$  and current avail- 522  
 able TxOp  $\mathbf{y}^t$ , independent of the states, TxOps, and actions 523  
 prior to the time  $t$ , i.e.,  $\pi_i^t(\mathcal{O}_i^t) = \pi_i^t(s_i^t, \mathbf{y}^t)$ . The policy  $\boldsymbol{\pi}_i$  524  
 is said to be stationary if the bidding policy  $\pi_i^t = \pi_i$  for 525  
 $\forall t$ . The reward  $r_i(s^k, \mathbf{y}^k, \mathbf{a}^k)$  of the stage  $k$  is discounted 526  
 by the factor  $(\alpha_i)^{k-t}$  at time  $t$ . The factor  $\alpha_i (0 \leq \alpha_i < 1)$  527  
 is the discounted factor determined by a specific application 528  
 (for instance, for video streaming applications, this factor can 529  
 be set based on the tolerable delay). The total discounted sum 530  
 of rewards  $Q_i^t(s^t, \mathbf{y}^t, \boldsymbol{\pi})$  for SU  $i$  can be calculated at time 531  
 $t$  starting from the state profile  $s^t$ , assuming that all SUs 532  
 deploy stationary and Markov policies  $\boldsymbol{\pi} = (\boldsymbol{\pi}_i, \boldsymbol{\pi}_{-i})$ , as in (7), 533  
 shown at the bottom of the next page. The total discounted 534  
 sum of rewards in (7) consists of two parts: 1) the current 535  
 stage reward and 2) the expected future reward discounted by 536  
 $\alpha_i$ . Note that SU  $i$  cannot independently determine the above 537  
 value without explicitly knowing the policies and states of other 538  
 SUs. The SU maximizes the total discounted sum of future 539  
 rewards to select the bidding policy, which explicitly considers 540

$$p_i^{\text{buf}}(v_i^{t+1} | v_i^t, z_i^t) = \begin{cases} \frac{(\mu_i \Delta T)^h e^{-\mu_i \Delta T}}{h!}, & \text{if } 0 \leq h < X_i - (v_i^t - n(s_i^t, z_i^t))^+ \\ \sum_{k=h}^{\infty} \frac{(\mu_i \Delta T)^k e^{-\mu_i \Delta T}}{k!}, & \text{if } h = X_i - (v_i^t - n(s_i^t, z_i^t))^+ \end{cases} \quad (5)$$

541 the impact of the current bid vector on the expected future  
542 rewards. We define the *best response*  $\beta_i$  for SU  $i$  to other SUs'  
543 policies  $\pi_{-i}$  as

$$\beta_i(\pi_{-i}) = \arg \max_{\pi_i} Q_i^t(\mathbf{s}^t, \mathbf{y}^t, (\pi_i, \pi_{-i})). \quad (8)$$

544 The central issue in our stochastic game is how the best-  
545 response policies can be determined by the SUs. In the repeated  
546 auction mechanism discussed in Section III, the procedure that  
547 each SU  $i$  follows to compete for the channel opportunities is  
548 illustrated in Fig. 3. In this procedure, the bidding strategy  $\pi_i^t$  is  
549 continuously improved by the ‘‘bidding strategy improvement’’  
550 module. In Section V-B, we discuss the challenges involved in  
551 building such a module, and in Section VI, we develop a best-  
552 response learning algorithm that can be used to improve the  
553 bidding strategy.

#### 554 B. Challenges for Selecting the Best-Response 555 Bidding Policy

556 Recall that during each time slot, the CSM announces an  
557 auction based on the available TxOps, and then SUs bid for  
558 the resources. To enable the successful deployment of this  
559 resource auction mechanism, we can prove (similar to our  
560 prior work in [21]) that SUs have no incentive to misrepresent  
561 their information, i.e., they adhere to the ‘‘truth telling’’ policy.  
562 We assume that at each time slot  $t$ , SU  $i$  has preference  $u_{ij}^t$   
563 over the channel  $j$ , which captures the benefit derived when  
564 using that channel. The preference  $u_{ij}^t$  is interpreted as the  
565 benefit obtained by SU  $i$  when using channel  $j$  compared to the  
566 benefit when this channel is not used. Note that this benefit also  
567 includes the expected future rewards. The optimal bid  $a_{ij}^{t,\text{opt}}$   
568 that SU  $i$  can take on channel  $j$  at time  $t$  is the bid maximizing  
569 the net benefit  $u_{ij}^t + \tau_i^t$ . In the auction discussed in Section III,  
570 the optimal bid that SU  $i$  can make is  $a_{ij}^{t,\text{opt}} = u_{ij}^t$ , i.e., the  
571 optimal bid for SU  $i$  is to announce its true preference to the  
572 CSM [21]. The proof is omitted here due to space limitations,  
573 since it is similar to that in [21]. The payment made by SU  $i$  is

computed by the CSM based on the inconvenience incurred by 574  
other SUs due to SU  $i$  during that time slot [23]. 575

Next, we define the preference  $u_{ij}^t$  in the context of the 576  
stochastic game model. Using the channel  $j$ , SU  $i$  obtains 577  
the immediate gain  $g_i^t(s_i^t, \mathbf{y}^t, e_j)$  by transmitting the pack- 578  
ets in its buffer, where  $e_j$  indicates that channel  $j$  is al- 579  
located to SU  $i$  during the current time slot. SU  $i$  then 580  
moves into the next state  $\mathbf{s}_i^{t+1}$  from which it may ob- 581  
tain the future reward  $Q_i^{t+1}(\mathbf{s}^{t+1}, \mathbf{y}^{t+1}, \pi)$ . On the other 582  
hand, if no channel is assigned to SU  $i$ , it receives the 583  
immediate gain  $g_i^t(s_i^t, \mathbf{y}^t, \mathbf{0})$  and then moves into the next 584  
state  $\mathbf{s}_i^{t+1}$ , from which it may obtain the future reward 585  
 $Q_i^{t+1}(\mathbf{s}^{t+1}, \mathbf{y}^{t+1}, \pi)$ . We define a feasible set of channel as- 586  
signments to SU  $i$ 's opponents (given SU  $i$ 's channel allocation 587  
 $\mathbf{z}_i^t$ ) as  $\mathcal{Z}_{-i}^t(\mathbf{z}_i^t)$ , with  $\mathcal{Z}_{-i}^t(\mathbf{z}_i^t) = \{Z_{-i}^t | \sum_{k=1, k \neq i}^M z_{kj}^t = y_j^t - 588$   
 $z_i^t \forall j, \sum_{j=1}^N z_{kj}^t \leq 1 \forall k \neq i, z_{kj}^t \in \{0, 1\}\}$ . 589

The preference over the current state can then be computed as 590

$$\begin{aligned} u_{ij}^t(\mathbf{s}^t, \mathbf{y}^t) &= \left[ g_i^t(s_i^t, \mathbf{y}^t, e_j) + \alpha_i \sum_{\substack{\mathbf{s}^{t+1} \in \mathcal{S} \\ \mathbf{y}^{t+1} \in \{0,1\}^N}} \right. \\ &\quad \times \left[ q_i(s_i^{t+1}, \mathbf{y}^{t+1} | s_i^t, \mathbf{y}^t, e_j) \sum_{Z_{-i}^t \in \mathcal{Z}_{-i}^t(e_j)} \right. \\ &\quad \times \left. \left. \left[ \prod_{k=1}^M q_k(s_k^{t+1}, \mathbf{y}^{t+1} | s_k^t, \mathbf{y}^t, z_k^t) Q_i^{t+1}(\mathbf{s}^{t+1}, \mathbf{y}^{t+1}, \pi) \right] \right] \right] \\ &\quad - \left[ g_i^t(s_i^t, \mathbf{y}^t, \mathbf{0}) + \alpha_i \sum_{\substack{\mathbf{s}^{t+1} \in \mathcal{S} \\ \mathbf{y}^{t+1} \in \{0,1\}^N}} \right. \\ &\quad \times \left[ q_i(s_i^{t+1}, \mathbf{y}^{t+1} | s_i^t, \mathbf{y}^t, \mathbf{0}) \sum_{Z_{-i}^t \in \mathcal{Z}_{-i}^t(\mathbf{0})} \right. \\ &\quad \times \left. \left. \left[ \prod_{k=1}^M q_k(s_k^{t+1}, \mathbf{y}^{t+1} | s_k^t, \mathbf{y}^t, z_k^t) Q_i^{t+1}(\mathbf{s}^{t+1}, \mathbf{y}^{t+1}, \pi) \right] \right] \right]. \end{aligned} \quad (9)$$

$$\begin{aligned} Q_i^t(\mathbf{s}^t, \mathbf{y}^t, \pi) &= \sum_{k=t}^{\infty} (\alpha_i)^{k-t} r_i(\mathbf{s}^k, \mathbf{y}^k, \pi(\mathbf{s}^k, \mathbf{y}^k)) = \underbrace{r_i(\mathbf{s}^t, \mathbf{y}^t, \pi(\mathbf{s}^t, \mathbf{y}^t))}_{\text{stage reward at time } t} \\ &\quad + \alpha_i \underbrace{\sum_{\substack{\mathbf{s}^{t+1} \in \mathcal{S} \\ \mathbf{y}^{t+1} \in \{0,1\}^N}} \left\{ \prod_{k=1}^M q_k(s_k^{t+1}, \mathbf{y}^{t+1} | s_k^t, \mathbf{y}^t, z_k^t) (\pi(\mathbf{s}^t, \mathbf{y}^t)) \times Q_i^{t+1}(\mathbf{s}^{t+1}, \mathbf{y}^{t+1}, \pi) \right\}}_{\text{expected future reward}} \\ &= \underbrace{\left\{ g_i^t(s_i^t, \mathbf{y}^t, z_i^t(\pi(\mathbf{s}^t, \mathbf{y}^t))) + \tau_i^t(\pi(\mathbf{s}^t, \mathbf{y}^t)) \right\}}_{\text{stage reward at time } t} \\ &\quad + \alpha_i \underbrace{\sum_{\substack{\mathbf{s}^{t+1} \in \mathcal{S} \\ \mathbf{y}^{t+1} \in \{0,1\}^N}} \left\{ \prod_{k=1}^M q_k(s_k^{t+1}, \mathbf{y}^{t+1} | s_k^t, \mathbf{y}^t, z_k^t) (\pi(\mathbf{s}^t, \mathbf{y}^t)) \times Q_i^{t+1}(\mathbf{s}^{t+1}, \mathbf{y}^{t+1}, \pi) \right\}}_{\text{expected future reward}} \end{aligned} \quad (7)$$

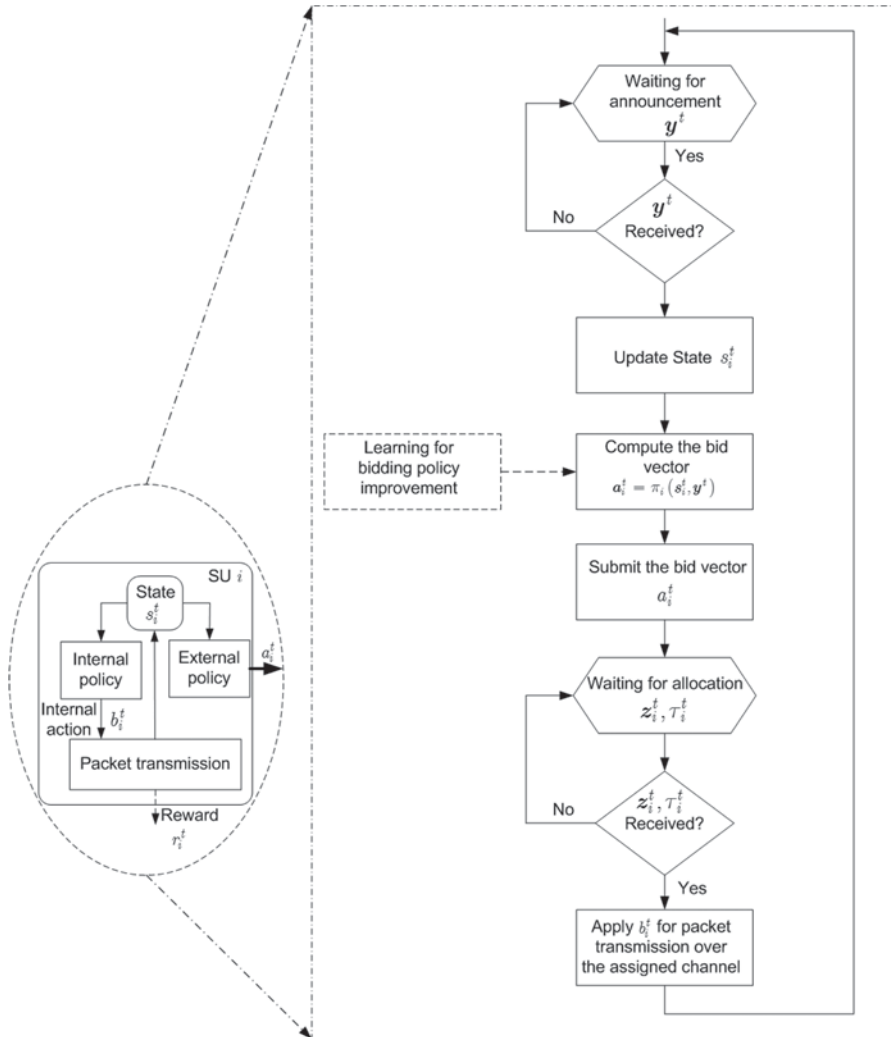


Fig. 3. Procedure for SU  $i$  to play the auction game at time slot  $t$ .

591 From this equation, it is clear that the true value  $u_{ij}^t$  depends  
 592 not only on its own current state  $s_i^t$  but also on the other SUs'  
 593 states  $s_{-i}^t$ , the channel allocations  $\mathcal{Z}_{-i}^t(e_j)$  to the other users  
 594 when channel  $j$  is assigned to SU  $i$ ,  $\mathcal{Z}_{-i}^t(\mathbf{0})$  when SU  $i$  is  
 595 not assigned to any channel, and the state transition models  
 596  $q_k(s_k^{t+1}, \mathbf{y}^{t+1} | s_k^t, \mathbf{y}^t, \mathbf{z}_k^t) \forall k$ . However, the other SUs' states,  
 597 the channel allocations, and the state transition models of other  
 598 SUs are not known to SU  $i$ , and it is, thus, impossible for each  
 599 SU to determine its preference  $u_{ij}^t(s^t, \mathbf{y}^t)$ .

600 Without knowing the other SUs' states and state transition  
 601 models, SU  $i$  cannot derive its optimal bidding strategy  
 602  $a_{ij}^{t, \text{opt}} = u_{ij}^t(s^t, \mathbf{y}^t)$ . However, if SU  $i$  chooses the bid  
 603 vector by only maximizing the immediate reward  $g_i^t + \tau_i^t$ ,  
 604 i.e., the total discounted sum of reward degenerates in  
 605  $Q_i^t(s^t, \mathbf{y}^t, \pi) = g_i^t(s_i^t, \mathbf{y}^t, \mathbf{z}_i^t(\pi(s^t, \mathbf{y}^t))) + \tau_i^t(\pi(s^t, \mathbf{y}^t))$  by  
 606 setting  $\alpha_i = 0$ . Then, the preference over channel  $j$  becomes  
 607  $u_{ij}^t(s^t, \mathbf{y}^t) = g_i^t(s_i^t, \mathbf{y}^t, e_j) - g_i^t(s_i^t, \mathbf{y}^t, \mathbf{0})$ . Now, since  $u_{ij}^t$   
 608 only depends on the state  $s_i^t$ , SU  $i$  can compute both the  
 609 optimal bid vector and the optimal bidding policy. We refer to  
 610 this optimal bidding policy as the "myopic" policy since it only  
 611 takes the immediate reward into consideration and ignores the  
 612 future impact. The myopic policy is referred to as  $\pi_i^{\text{myopic}}$ . To

solve the difficult problem of optimal bidding policy selection  
 613 when  $\alpha_i \neq 0$ , an SU needs to forecast the impact of its current  
 614 bidding actions on the expected future rewards discounted by  
 615  $\alpha_i$ . The forecast can be performed using learning from its past  
 616 experiences. 617

## VI. INTERACTIVE LEARNING FOR PLAYING THE RESOURCE MANAGEMENT GAME

### A. How to Evaluate Learning Algorithms?

621 Section V-B shows that an SU needs to know the other SUs'  
 622 states and state transition models to derive its own optimal  
 623 bidding policy. This coupling among SUs is due to the shared  
 624 nature of wireless resources. However, an SU cannot exactly  
 625 know the other SUs' models and private information in wireless  
 626 networks. Thus, to improve the bidding policy, an SU can only  
 627 predict the impacts of dynamics (uncertainties) caused by the  
 628 competing SUs based on its observations from past auctions. 628  
 In this paper, we propose a learning algorithm for predicting  
 629 these impacts. We define a learning algorithm  $\mathcal{L}_i$  for SU  $i$  as  
 630 a function taking the observation  $\mathcal{O}_i^t$  as input and having the  
 631 bidding policy  $\pi_i^t$  as output. 632



633 Before developing a learning algorithm, we first discuss how  
634 to evaluate the performance of a learning algorithm in terms  
635 of its impact on the SU's reward. Unlike existing multiagent  
636 learning research, which is aimed at achieving converge to an  
637 equilibrium point for the interacting agents, we develop learn-  
638 ing algorithms based on the performance of the bidding strategy  
639 on the SU's reward. We denote a bidding policy generated by  
640 the learning algorithm  $\mathcal{L}_i$  as  $\pi_i^{\mathcal{L}_i}$ . An SU will learn to improve  
641 its bidding policy and its rewards from participating in the  
642 auction game. The performance of the bidding strategy  $\pi_i$  is  
643 defined as the time average reward that SU  $i$  obtains in a time  
644 window with length  $T$  when it adopts  $\pi_i$ , i.e.,

$$\mathcal{V}^{\pi_i}(T) = \frac{1}{T} \sum_{k=1}^T r_i^k. \quad (10)$$

645 Using this definition, the performance of two learning al-  
646 gorithms can easily be compared. For instance, given two  
647 algorithms  $\mathcal{L}'_i$  and  $\mathcal{L}''_i$ , if  $\mathcal{V}^{\pi_i^{\mathcal{L}'_i}} > \mathcal{V}^{\pi_i^{\mathcal{L}''_i}}$ , then we say that the  
648 learning algorithm  $\mathcal{L}'_i$  is better than  $\mathcal{L}''_i$ .

#### 649 B. What Information to Learn From?

650 First, let us consider what information the SU can  
651 observe while playing the stochastic game in our SN. As  
652 shown in Fig. 1, at the beginning of time slot  $t$ , the SUs  
653 submit the bids  $a_i^t \forall i$ . Then, the CSM returns the channel  
654 allocations  $z_i^t \forall i$  and  $\tau_i^t \forall i$ . If SU  $i$  is not allowed to  
655 observe the bids, the channel allocations, and payments  
656 for other SUs, then the observation of SU  $i$  becomes  $\mathbf{o}_i^t =$   
657  $\{\mathbf{s}_i^0, \mathbf{y}^0, a_i^0, z_i^0, \tau_i^0, \dots, \mathbf{s}_i^{t-1}, \mathbf{y}^{t-1}, a_i^{t-1}, z_i^{t-1}, \tau_i^{t-1}, \mathbf{s}_i^t, \mathbf{y}^t\}$ . If  
658 the information is exchanged among SUs or broad-  
659 casted and overheard by all SUs, the observed infor-  
660 mation by SU  $i$  becomes  $\mathbf{o}_i^t = \{\mathbf{s}_i^0, \mathbf{y}^0, a_i^0, z_i^0, \tau_i^0, \dots,$   
661  $\mathbf{s}_i^{t-1}, \mathbf{y}^{t-1}, a_i^{t-1}, z_i^{t-1}, \tau_i^{t-1}, \mathbf{s}_i^t, \mathbf{y}^t\}$ . Now, the problem that  
662 needs to be solved by SU  $i$  is how it can improve its own policy  
663 for playing the game by learning from the observation  $\mathbf{o}_i^t$ . In  
664 this paper, we assume that SU  $i$  observes the information  $\mathbf{o}_i^t =$   
665  $\{\mathbf{s}_i^0, \mathbf{y}^0, a_i^0, z_i^0, \tau_i^0, \dots, \mathbf{s}_i^{t-1}, \mathbf{y}^{t-1}, a_i^{t-1}, z_i^{t-1}, \tau_i^{t-1}, \mathbf{s}_i^t, \mathbf{y}^t\}$ .

#### 666 C. What to Learn?

667 In Section VI-A, we introduce learning as a tool to predict the  
668 impacts of dynamics and, hence, improve the bidding policy.  
669 However, a key question is what needs to be learned. Recall that  
670 the optimal bidding policy for SU  $i$  is to generate a bid vector  
671 that represents its preferences for using different channels.  
672 From (9), we can see that SU  $i$  needs to learn the following:  
673 1) the state space of other SUs, i.e.,  $\mathcal{S}_{-i}$ ; 2) the current state of  
674 other SUs, i.e.,  $\mathbf{s}_{-i}^t$ ; 3) the transition probability of other SUs,  
675 i.e.,  $\prod_{k \neq i} q_k(\mathbf{s}_k^{t+1}, \mathbf{y}^{t+1} | \mathbf{s}_k^t, \mathbf{y}^t, \mathbf{z}_k^t)$ ; 4) the resource allocations  
676  $\mathcal{Z}_{-i}^t(e_j) \forall j$  and  $\mathcal{Z}_{-i}^t(\mathbf{0})$ ; and 5) the discounted sum of rewards  
677  $Q_i^{t+1}(\mathbf{s}^{t+1}, \mathbf{y}^{t+1}, \boldsymbol{\pi})$ .

678 However, SU  $i$  can only observe the information  $\mathbf{o}_i^t = \{\mathbf{s}_i^0,$   
679  $\mathbf{y}^0, a_i^0, z_i^0, \tau_i^0, \dots, \mathbf{s}_i^{t-1}, \mathbf{y}^{t-1}, a_i^{t-1}, z_i^{t-1}, \tau_i^{t-1}, \mathbf{s}_i^t, \mathbf{y}^t\}$  from  
680 which SU  $i$  cannot accurately infer the other SUs' state space  
681 and transition probability. Moreover, capturing the exact in-

formation about other SUs requires heavy computational and 682  
storage complexity. Instead, we allow SU  $i$  to classify the space 683  
 $\mathcal{S}_{-i}$  into  $H_i$  classes, each of which is represented by a represen- 684  
tative state  $\tilde{\mathbf{s}}_{-i,h}, h \in \{1, \dots, H_i\}$ . We discuss how the space 685  
 $\mathcal{S}_{-i}$  is decomposed in Section VI-D. By dividing the state space 686  
 $\mathcal{S}_{-i}$ , the transition probability  $\prod_{k \neq i} q_k(\mathbf{s}_k^{t+1}, \mathbf{y}^{t+1} | \mathbf{s}_k^t, \mathbf{y}^t, \mathbf{z}_k^t)$  687  
is approximated by  $q_{-i}(\tilde{\mathbf{s}}_{-i}^{t+1}, \mathbf{y}^{t+1} | \tilde{\mathbf{s}}_{-i}^t, \mathbf{y}^t, \mathbf{z}_i^t)$ , where  $\tilde{\mathbf{s}}_{-i}^t$  and 688  
 $\tilde{\mathbf{s}}_{-i}^{t+1}$  are the representative states of the classes to which  $\mathbf{s}_{-i}^t$  and 689  
 $\mathbf{s}_{-i}^{t+1}$  belong. This approximation is performed by aggregating 690  
all the other SUs' states into one representative state and assum- 691  
ing that the transition depends on the resource allocation  $\mathbf{z}_i^t$ . 692  
The transition probability approximation is also discussed in 693  
Section VI-D. The discounted sum of rewards  $Q_i^{t+1}(\mathbf{s}^{t+1},$  694  
 $\mathbf{y}^{t+1}, \boldsymbol{\pi})$  is approximated by  $V_i^{t+1}((\mathbf{s}_i^{t+1}, \tilde{\mathbf{s}}_{-i}^{t+1}), \mathbf{y}^{t+1})$ . 695  
Note that the classification on the state space  $\mathcal{S}_{-i}$  and the 696  
approximation of the transition probability and discounted sum 697  
of rewards affect the learning performance. Hence, a user can 698  
tradeoff an increased complexity for an increased performance. 699  
After the classification, the preference computation can be 700  
approximated as 701

$$\begin{aligned} & u_{ij}^t((\mathbf{s}_i^t, \tilde{\mathbf{s}}_{-i}^t), \mathbf{y}^t) \\ &= \left[ g_i^t q(\mathbf{s}_i^t, \mathbf{y}^t, e_j) + \alpha_i \sum_{\substack{(\mathbf{s}_i^{t+1}, \tilde{\mathbf{s}}_{-i}^{t+1}) \in (\mathcal{S}_i, \tilde{\mathcal{S}}_{-i}) \\ \mathbf{y}^{t+1} \in \{0,1\}^N}} \right. \\ & \quad \times \left[ q_i(\mathbf{s}_i^{t+1}, \mathbf{y}^{t+1} | \mathbf{s}_i^t, \mathbf{y}^t, e_j) \times q_{-i}(\tilde{\mathbf{s}}_{-i}^{t+1}, \mathbf{y}^{t+1} | \tilde{\mathbf{s}}_{-i}^t, \mathbf{y}^t, e_j) \right. \\ & \quad \left. \left. \times V_i^{t+1}((\mathbf{s}_i^{t+1}, \tilde{\mathbf{s}}_{-i}^{t+1}), \mathbf{y}^{t+1}) \right] \right] \\ & - \left[ g_i^t(\mathbf{s}_i^t, \mathbf{y}^t, \mathbf{0}) + \alpha_i \sum_{\substack{(\mathbf{s}_i^{t+1}, \tilde{\mathbf{s}}_{-i}^{t+1}) \in (\mathcal{S}_i, \tilde{\mathcal{S}}_{-i}) \\ \mathbf{y}^{t+1} \in \{0,1\}^N}} \right. \\ & \quad \times \left[ q_i(\tilde{\mathbf{s}}_{-i}^{t+1}, \mathbf{y}^{t+1} | \mathbf{s}_i^t, \mathbf{y}^t, \mathbf{0}) \times q_{-i}(\tilde{\mathbf{s}}_{-i}^{t+1}, \mathbf{y}^{t+1} | \tilde{\mathbf{s}}_{-i}^t, \mathbf{y}^t, \mathbf{0}) \right. \\ & \quad \left. \left. \times V_i^{t+1}((\mathbf{s}_i^{t+1}, \tilde{\mathbf{s}}_{-i}^{t+1}), \mathbf{y}^{t+1}) \right] \right]. \quad (11) \end{aligned}$$

In this setting, to find the approximated preference and, 702  
thus, the approximated optimal bidding policy, we need 703  
to learn the following from past observations: 1) how 704  
the space  $\tilde{\mathcal{S}}_{-i}$  is classified; 2) the transition probability 705  
 $q_{-i}(\tilde{\mathbf{s}}_{-i}^{t+1}, \mathbf{y}^{t+1} | \tilde{\mathbf{s}}_{-i}^t, \mathbf{y}^t, \mathbf{z}_i^t)$ ; and 3) the approximated future 706  
rewards  $V_i^{t+1}((\mathbf{s}_i^{t+1}, \tilde{\mathbf{s}}_{-i}^{t+1}), \mathbf{y}^{t+1})$ . 707

#### 708 D. How to Learn?

In this section, we develop a learning algorithm to estimate 709  
the terms listed in Section VI-C. 710

1) *Decomposition of the Space  $\mathcal{S}_{-i}$* : As discussed 711  
in Section VI-B, only  $\mathbf{o}_i^t = \{\mathbf{s}_i^0, \mathbf{y}^0, a_i^0, z_i^0, \tau_i^0, \dots, \mathbf{s}_i^{t-1},$  712  
 $\mathbf{y}^{t-1}, a_i^{t-1}, z_i^{t-1}, \tau_i^{t-1}, \mathbf{s}_i^t, \mathbf{y}^t\}$  are observed. From the auction 713  
mechanism presented in Section III, we know that the value of 714

715 the tax  $\tau_i^t$  is computed based on the inconvenience that SU  $i$   
 716 causes to the other SUs. In other words, a higher value of  $|\tau_i^t|$   
 717 indicates that the network is more congested.<sup>7</sup> Based on the  
 718 bid vector  $\mathbf{b}_i^t$ , the channel allocation  $\mathbf{z}_i^t$ , and the tax  $\tau_i^t$ , SU  $i$   
 719 can infer network congestion and thus, indirectly, the resource  
 720 requirements of the competing SUs. Instead of knowing the  
 721 exact state space of other SUs, SU  $i$  can classify the space  $\mathcal{S}_{-i}$   
 722 as follows.

723 We assume that the maximum absolute tax is  $\Gamma$ . We split the  
 724 range  $[0, \Gamma]$  into  $[\Gamma_0, \Gamma_1), [\Gamma_1, \Gamma_2), \dots, [\Gamma_{H_i-1}, \Gamma_{H_i}]$  with  $0 =$   
 725  $\Gamma_0 \leq \Gamma_1 \leq \dots \leq \Gamma_{H_i} = \Gamma$ . Here, we assume that the values  
 726 of  $\{\Gamma_1, \dots, \Gamma_{H_i-1}\}$  are equally located in the range of  $[0, \Gamma]$ .  
 727 (Note that more sophisticated selection for these values can be  
 728 deployed, and this forms an interesting area of future research.)

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

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

$$\tilde{s}_{-i}^t = h, \quad \text{if } |\tau_i^t| \in [\Gamma_{h-1}, \Gamma_h). \quad (12)$$

733 2) If the resource allocation  $\mathbf{z}_i^t = \mathbf{0}$  but  $\mathbf{y}^t \neq \mathbf{0}$ , the tax is  
 734 0. In this case, we cannot use the tax to predict network  
 735 congestion. However, we can infer that the congestion  
 736 is more severe than the minimum bid for those avail-  
 737 able channels, i.e.,  $\min_{j \in \{l: y_l^t \neq 0\}} \{a_{ij}^t\}$ . This is because,  
 738 in this current stage of the auction game, only SU  $i'$   
 739 with  $a_{i'j}^t \geq a_{ij}^t$  can obtain channel  $j$ , which indicates  
 740 that  $|\tau_{i'}^t| \geq \min_{j \in \{l: y_l^t \neq 0\}} \{a_{ij}^t\}$  if SU  $i$  is allocated any  
 741 channel. Then, the representative state of the other SUs  
 742 is chosen as

$$\tilde{s}_{-i}^t = h, \quad \text{if } \min_{j \in \{l: y_l^t \neq 0\}} \{a_{ij}^t\} \in [\Gamma_{h-1}, \Gamma_h). \quad (13)$$

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

<sup>7</sup>When the CSM deploys a mechanism without tax for resource management, the space classification for other SUs can also be done based on the announced information and corresponding resource allocation.

2) *Estimating the Transition Probability:* To estimate the 746  
 transition probability, SU  $i$  maintains a table  $F$  with size  $H_i \times 747$   
 $H_i \times (N + 1)$ . Each entry  $f_{h', h'', j}$  in the table  $F$  represents the 748  
 number of transitions from state  $\tilde{s}_{-i}^t = h''$  to state  $\tilde{s}_{-i}^{t+1} = h'$  749  
 when the resource allocation  $\mathbf{z}_i^t = \mathbf{e}_j$  (or  $\mathbf{0}$  if  $j = 0$ ). It is 750  
 clear that  $H_i$  will significantly influence the complexity and 751  
 memory requirements, etc., of SU  $i$ . The update of  $F$  is simply 752  
 based on the observation  $\sigma_i^t$  and the state classification in the 753  
 above section. Then, we use the frequency to approximate the 754  
 transition probability [15], i.e., 755

$$q_{-i}(\tilde{s}_{-i}^{t+1} = h' | \tilde{s}_{-i}^t = h'', \mathbf{e}_j) = \frac{f_{h', h'', j}}{\sum_{h'} f_{h', h'', j}}. \quad (14)$$

3) *Learning the Future Reward:* By classifying the state 756  
 space  $\mathcal{S}_{-i}$  and estimating the transition probability, SU  $i$  757  
 can now forecast the value of the average future reward 758  
 $V_i^{t+1}((s_i^{t+1}, \tilde{s}_{-i}^{t+1}), \mathbf{y}^{t+1})$  using learning. Equation (7) can be 759  
 approximated by (15), shown at the bottom of the page. 760

Similar to the Q-learning established in [17], we also use 761  
 the received rewards to update the estimation of future rewards. 762  
 However, the main difference between our proposed algorithm 763  
 and Q-learning is that our solution explicitly considers the 764  
 impacts of other SUs' bidding actions through the state clas- 765  
 sifications and transition probability approximation. 766

We use a 3-D table to store the value  $V_i((s_i, \tilde{s}_{-i}), \mathbf{y})$  with 767  
 $s_i \in \mathcal{S}_i$ ,  $\tilde{s}_{-i} \in \tilde{\mathcal{S}}_{-i}$ . The total number of entries in  $V_i$  is  $|\mathcal{S}_i| \times 768$   
 $H_i \times 2^N$ . SU  $i$  updates the value of  $V_i((s_i, \tilde{s}_{-i}), \mathbf{y})$  at time 769  
 $t$  according to the rules in (16), shown at the bottom of the 770  
 page, where  $\gamma_i^t \in [0, 1)$  is a learning rate factor satisfying 771  
 $\sum_{t=1}^{\infty} \gamma_i^t = \infty$ , and  $\sum_{t=1}^{\infty} (\gamma_i^t)^2 < \infty$  [17]. In summary, the 772  
 learning procedure that is developed for an SU is shown in 773  
 Table I. 774

### E. Complexity of Learning 775

In Section III, we have discussed the computation complexity 776  
 incurred by the CSM and the communication cost between 777  
 the CSM and the SUs. In this section, we further quantify 778  
 the complexity of learning in terms of the computational and 779  
 storage burden. We use a floating-point operation ("flop") as a 780  
 measure of complexity, which will provide us an estimation of 781

$$Q_i^t((s_i^t, \tilde{s}_{-i}^t), \mathbf{y}^t, \boldsymbol{\pi}) \doteq \left\{ g_i^t(s_i^t, \mathbf{y}^t, \mathbf{z}_i^t(\boldsymbol{\pi}(s^t, \mathbf{y}^t))) + \tau_i^t(\boldsymbol{\pi}(s^t, \mathbf{y}^t)) + \alpha_i \sum_{\substack{(s_i^{t+1}, \tilde{s}_{-i}^{t+1}) \in (\mathcal{S}_i, \tilde{\mathcal{S}}_{-i}) \\ \mathbf{y}^{t+1} \in \{0, 1\}^N}} \right. \\ \left. \times \left\{ q_i(s_i^{t+1}, \mathbf{y}^{t+1} | s_i^t, \mathbf{y}^t, \mathbf{z}_i^t(\boldsymbol{\pi}(s^t, \mathbf{y}^t))) q_{-i}(\tilde{s}_{-i}^{t+1}, \mathbf{y}^{t+1} | \tilde{s}_{-i}^t, \mathbf{y}^t, \mathbf{z}_i^t(\boldsymbol{\pi}(s^t, \mathbf{y}^t))) V_i^{t+1}((s_i^{t+1}, \tilde{s}_{-i}^{t+1}), \mathbf{y}^{t+1}) \right\} \right\} \quad (15)$$

$$V_i^t((s_i, \tilde{s}_{-i}), \mathbf{y}) = \begin{cases} (1 - \gamma_i^t) V_i^{t-1}((s_i, \tilde{s}_{-i}), \mathbf{y}) + \gamma_i^t Q_i^t((s_i, \tilde{s}_{-i}), \mathbf{y}, \boldsymbol{\pi}), & \text{if } (s_i^t, \tilde{s}_{-i}^t) = (s_i, \tilde{s}_{-i}), \quad \mathbf{y}^t = \mathbf{y} \\ V_i^{t-1}((s_i, \tilde{s}_{-i}), \mathbf{y}), & \text{otherwise} \end{cases} \quad (16)$$

TABLE I  
LEARNING PROCEDURE

<p><b>Initializing:</b> <math>V_i^0((s_i, \tilde{s}_{-i}), \mathbf{y}) \leftarrow 0</math> for all possible states <math>s_i \in S_i, \tilde{s}_{-i} \in \tilde{S}_{-i}</math>.</p> <p><b>Learning:</b></p> <p>At time <math>t</math>, SU <math>i</math>:</p> <ol style="list-style-type: none"> <li>Observes the current state <math>s_i^t</math> and <math>\mathbf{y}^t</math>;</li> <li>Chooses an action <math>a_i^t = [u_{i1}^t, \dots, u_{iN}^t]</math> as computed in Eq. (11) by replacing <math>V_i^{t+1}((s_i^{t+1}, \tilde{s}_{-i}^{t+1}), \mathbf{y}^{t+1})</math> with <math>V_i^{t-1}((s_i^{t+1}, \tilde{s}_{-i}^{t+1}), \mathbf{y}^{t+1})</math>, and then submits it to the CSM;</li> <li>Receives the allocation <math>\mathbf{z}_i^t</math> and payment <math>\tau_i^t</math>;</li> <li>Computes the representative state <math>\tilde{s}_{-i}^t</math> as in Section VI.D.1) and update the transition probability as in Section VI.D.2);</li> <li>Computes the expected total discounted sum of the rewards <math>Q_i^t((s_i^t, \tilde{s}_{-i}^t), \mathbf{y}^t, \boldsymbol{\pi})</math> as in Eq. (15);</li> <li>Updates the future reward table <math>V_i^t((s_i, \tilde{s}_{-i}), \mathbf{y})</math> at the state <math>(s_i^t, \tilde{s}_{-i}^t)</math> and TxOp <math>\mathbf{y}^t</math> using the learning rate factor <math>\gamma_i^t</math>, according to Eq. (16).</li> </ol>
---

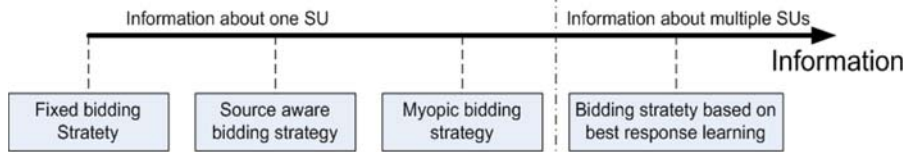


Fig. 4. Bidding strategies based on the required information.

782 the computational complexity required to perform the learning  
783 algorithm. In addition, based on this, we can determine how  
784 complexity grows with the increasing number of SUs [20]. At  
785 each stage, SU performs the classification of other the SUs'  
786 states, which, in the worst case, requires a number of "flops"  
787 approximately  $N$ . The number of "flops" to estimate the transi-  
788 tion probability of other SUs' states as in (14) is approximately  
789  $(H_i + 1)$ . The number of "flops" to learn the future reward  
790 is approximately  $(2|S_i|H_i + 6)$ . Therefore, the total number  
791 of "flops" incurred by the SU is  $N + H_i + 2|S_i|H_i + 7$ , from  
792 which we can note that the complexity of learning for each SU  
793 is proportional to the number of possible states of that SU and  
794 the number of classes in which the other SUs' state space is  
795 decomposed.

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

802

## VII. SIMULATION RESULTS

803 In this section, we aim at quantifying the performance of  
804 our proposed stochastic interaction and learning framework. We  
805 assume that the SUs compete for available spectrum opportuni-  
806 ties to transmit delay-sensitive multimedia data. First, we com-  
807 pare the performance of various bidding strategies. Next, we  
808 quantify the performance of our proposed learning algorithm  
809 in various network environments. We will only present here  
810 several illustrative examples. However, the same observations  
811 can be obtained using a larger number of SUs or channels.

### A. Various Bidding Strategies for Dynamic Multiuser Interaction

812

813

In this section, we highlight the merits of the stochastic  
814 game framework proposed in Section II by comparing the  
815 performance of different SUs, which deploy different bidding  
816 strategies. The SUs are required to submit the bid vector on  
817 the available channels. The SUs can deploy different bidding  
818 strategies to generate their bid vector. 819

- 1) Fixed bidding strategy  $\pi_i^{\text{fixed}}$ : This strategy generates a  
820 constant bid vector during each stage of the auction game,  
821 irrespective of the state that SU  $i$  is currently in and of the  
822 states other SUs are in. In other words,  $\pi_i^{\text{fixed}}$  does not  
823 consider any of the dynamics defined in Section IV. 824
- 2) Source-aware bidding strategy  $\pi_i^{\text{source}}$ : This strategy gen-  
825 erates various bid vectors by considering the dynamics in  
826 source characteristics (based on the current buffer state)  
827 but not the channel dynamics. 828
- 3) Myopic bidding strategy  $\pi_i^{\text{myopic}}$ : This strategy takes  
829 into account the disturbance due to the environment as  
830 well as the impact caused by other SUs, as discussed in  
831 Section V-B. However, it does not consider the impact on  
832 future rewards. 833
- 4) Bidding strategy based on best-response learning  $\pi_i^{\text{C}_i}$ :  
834 This strategy is produced using the learning algorithm  
835 proposed in Section VI.  $\pi_i^{\text{C}_i}$  considers the two types of  
836 dynamics defined in Section IV and the interaction impact  
837 on future reward. 838

In terms of required information, the above bidding strategies  
839 are illustrated in Fig. 4. For instance, the fixed bidding strategy  
840  $\pi_i^{\text{fixed}}$  does not require information about SU  $i$ 's state or other  
841 SUs' states. The source-aware bidding strategy  $\pi_i^{\text{buff}}$  considers 842

TABLE II  
PERFORMANCE OF SU 1 AND 2 WITH VARIOUS BIDDING STRATEGIES IN THE TWO SU NETWORKS

	Bidding Strategies	SU 1			SU 2		
		Packet loss rate (%)	Average tax	Average cost	Packet loss rate (10%)	Average tax	Average cost
Scenario 1	$\pi_1^{fixed}, \pi_2^{fixed}$	32.53	0.4875	2.8966	31.05	0.5095	2.6104
Scenario 2	$\pi_1^{fixed}, \pi_2^{myopic}$	34.36	0.1222	2.6337	14.39	0.5495	1.5105
Scenario 3	$\pi_1^{source}, \pi_2^{myopic}$	29.83	0.3147	2.4915	18.11	0.6048	1.6116
Scenario 4	$\pi_1^{myopic}, \pi_2^{myopic}$	21.55	0.4669	1.9767	19.55	0.3763	1.7837
Scenario 5	$\pi_1^{\mathcal{L}_1}, \pi_2^{myopic}$	15.14	0.6923	1.7428	27.29	0.4197	2.2967

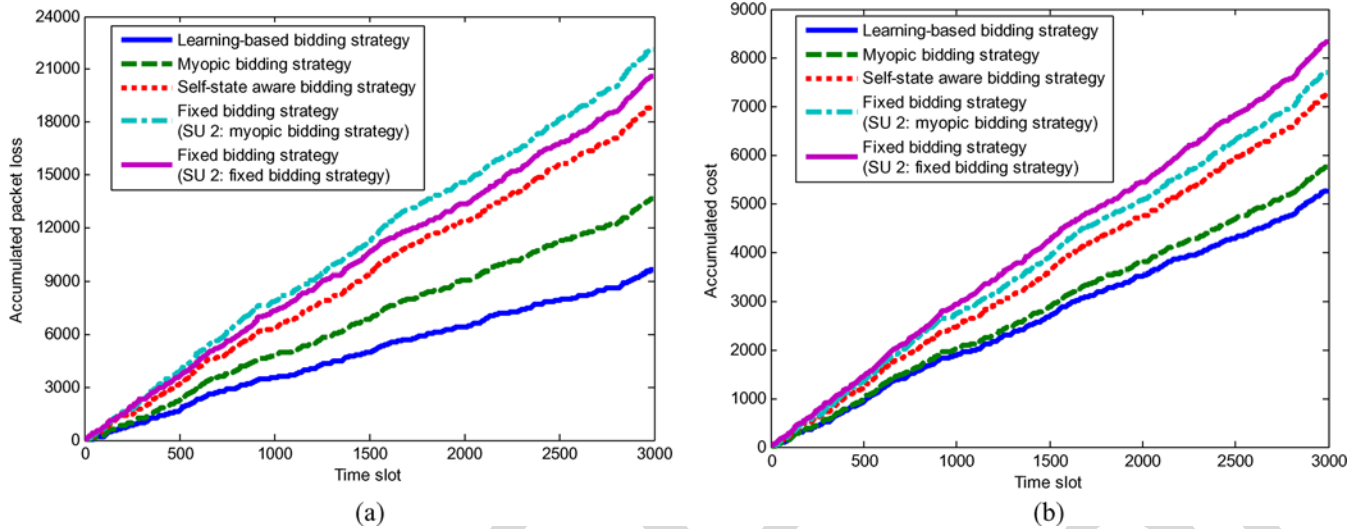


Fig. 5. Accumulated packet loss and cost of SU 1 in the five scenarios. (a) Accumulated packet loss over the time slot. (b) Accumulated cost over the time slot.

843 the source characteristics based on the current buffer state.  
844 However, the myopic bidding strategy  $\pi_i^{myopic}$  requires full  
845 information about SU  $i$ 's state. The bidding strategy based on  
846 best-response learning  $\pi_i^{\mathcal{L}_i}$  also requires information about the  
847 states of other SUs.

848 In this simulation, we consider the SN as an extension of  
849 WLANs with spectral agile capability [9]. In the following,  
850 we first simulate the case that two SUs compete for the chan-  
851 nel opportunities and then extend to the case with multiple  
852 (five) SUs.

853 1) *Competition Among Two SUs for Channel Opportunities:*

854 We first consider a simple illustrative network with two SUs  
855 competing for available TxOps. The packet arrivals of the SUs  
856 are modeled using a Poisson process with the same average  
857 arrival rate of 1 Mb/s. For simplicity of illustration, the channel  
858 condition of SU 1 (SU 2) on each channel only takes three val-  
859 ues ( $K = 3$ ), which are 18, 23, and 26 dB. The transition prob-  
860 abilities are  $p_{ij}^{0 \rightarrow 1} = p_{ij}^{0 \rightarrow 2} = 0.4$ ,  $p_{ij}^{0 \rightarrow 3} = 0.2$ ,  $p_{1j}^{1 \rightarrow 1} = p_{1j}^{1 \rightarrow 2} =$   
861  $0.4$ , and  $p_{1j}^{1 \rightarrow 3} = 0.2 \forall i, j, l$ . The transition probability of the  
862 availability of channels to SUs is  $p_j^{NF} = p_j^{FN} = 0.5$ . For sim-  
863 plicity of illustration, the environment parameters experienced  
864 by the two SUs are the same. The length of the time slot  $\Delta T$   
865 is  $10^{-2}$  s.

866 In this simulation, we consider five scenarios. In scenario 1,  
867 both SU 1 and SU 2 deploy the fixed bidding strategy  $\pi_1^{fixed}$ .  
868 In scenarios 2–5, SU 1 deploys the fixed bidding strategy  
869  $\pi_1^{fixed}$ , source-aware bidding strategy  $\pi_1^{source}$ , myopic bidding

strategy  $\pi_1^{myopic}$ , and best-response learning-based bidding  
870 strategy  $\pi_1^{\mathcal{L}_1}$ , respectively, and SU 2 always deploys the myopic  
871 bidding strategy  $\pi_2^{myopic}$ . The discounted factor for the best-  
872 response learning algorithm is set to 0.8. As discussed in  
873 Section IV-B, the stage reward is defined as  $r_i^t = (g_i^t + \tau_i^t)$ ,  
874 with  $(-g_i^t - \tau_i^t)$  being the number of packet lost plus the tax  
875 charged by the CSM (note that  $\tau_i^t \leq 0$ ). This can be interpreted  
876 as the cost incurred at each stage. Similar to (10), we use the  
877 average cost over the time window  $T = 1000$  to evaluate the  
878 performance of the bidding strategies. Hence, the lower  
879 the average cost, the better the performance of the bidding  
880 strategy. The packet loss rate, average tax, and cost per time slot  
881 are presented in Table II. The accumulated packet loss and cost  
882 of SU 1 for the five scenarios are plotted in Fig. 5(a) and (b),  
883 respectively.

884 From this simulation, comparing scenario 2 with scenario 1,  
885 we observe that when SU 2 deploys the myopic strategy against  
886 SU 1, which adopted the fixed bidding strategy, SU 2 reduces  
887 its average cost by around 42% and the average packet loss  
888 rate by around 16.6%. This significant improvement is because  
889 SU 2 can more accurately value the channel opportunities by  
890 modeling and considering its experienced dynamics, i.e., source  
891 characteristics, channel conditions, and availability.

892 In scenario 3, SU 1 improves its bidding strategy (i.e.,  
893 it deploys now a source-aware bidding strategy) by partially  
894 considering its experienced environment, i.e., SU 1 generates  
895 its bid vector by only considering the source dynamics though  
896

TABLE III  
PERFORMANCE OF SU 1–5 WITH VARIOUS BIDDING STRATEGIES IN THE FIVE SU NETWORKS

	SU 1		SU 2		SU 3		SU 4		SU 5	
	Packet Loss Rate (%)	Average cost	Packet Loss Rate (%)	Average cost	Packet Loss Rate (%)	Average cost	Packet Loss Rate (%)	Average cost	Packet Loss Rate (%)	Average cost
1	21.14	1.2002	19.99	1.1666	22.05	1.2123	21.37	1.1949	24.17	1.3101
2	25.03	1.2992	24.20	1.2993	25.72	1.3338	26.02	1.3568	9.56	1.0988

897 its current buffer state. Compared with scenario 2, if SU 1  
898 considers more information about its own state, it can further  
899 reduce its packet loss rate by an average of 4.5% and an  
900 average cost by around 5.4%. This observation verifies that the  
901 information about the SU's state improves the bidding strategy.

902 In scenario 4, SU 1 deploys a myopic bidding strategy, which  
903 is more advanced than the source-aware bidding strategy since  
904 it considers both types of dynamics defined in Section IV  
905 (including the dynamics regarding the source characteristics,  
906 channel conditions, and channel availability, and the interaction  
907 with other SUs in the auction mechanism). The significant  
908 improvement in terms of packet loss rate (13% reduced) and  
909 average cost (25% reduced), compared with scenario 2, indi-  
910 cates that the myopic bidding strategy provides the optimal bid  
911 vector when only current benefits are considered, as shown in  
912 Section V-B.

913 In scenario 5, SU 1 further improves the bidding strat-  
914 egy using the best-response learning algorithm developed in  
915 Section VI. Using learning, SU 1 reduces the packet loss rate to  
916 15.14% and the average cost to 1.7428 (11.8% lower compared  
917 with scenario 4). This significant improvement is due to the  
918 ability of the SU to learning and forecast the future impact of  
919 its current actions.

920 It is also worth noting that the reduction of packet loss rate  
921 of SU 1 in scenarios 2–5 comes from two parts: One is the  
922 advanced bidding strategies, which allows the SU to take into  
923 consideration more information about its own states and the  
924 other SUs' states and, based on this better forecast, the impact  
925 of various actions; the other one is the increase in the amount  
926 of resources consumed by SU 1, which corresponds to a higher  
927 tax charged by the CSM, as shown in Table II.

928 We further note that the bidding strategy deployed by SU 1  
929 will affect the performance of SU 2. For example, comparing  
930 scenario 2 with scenario 4, the fixed bidding strategy of SU 1  
931 in scenario 2 leads to a lower average cost (15% reduced) for  
932 SU 2. This is because SU 1 uses a fixed bidding strategy, which  
933 does not account for the dynamic changes in its environment,  
934 while SU 2 minimizes its current cost (the number of packets  
935 lost plus the tax) based on its current state. However, when  
936 comparing scenario 5 with scenario 4, SU 1 using learning  
937 not only improves its prediction of the current environment  
938 dynamics but also better predicts the impact on the future cost  
939 based on the observations. The improvement leads to higher  
940 resource allocation (hence, incurring higher tax, see in Table II)  
941 for SU 1, thereby resulting in worse performance for SU 2 (i.e.,  
942 the average cost is increased by 22.2%).

943 2) *Multiple SUs Competition for Channel Opportunities:*  
944 In this simulation, we consider five SUs competing for the  
945 available TxOps in the WLAN-like SN. The packet arrivals of

all the five SUs are modeled using a Poisson process with the 946  
same average arrival rate of 1 Mb/s. The number of channels 947  
is 3, and the channel condition of all the five SUs on each 948  
channel takes only three values ( $K = 3$ ), which are 18, 23, 949  
and 26 dB. The transition probabilities are  $p_{ij}^{0 \rightarrow 1} = p_{ij}^{0 \rightarrow 2} = 0.4$ , 950  
 $p_{ij}^{0 \rightarrow 3} = 0.2$ ,  $p_{1j}^{1 \rightarrow 1} = p_{1j}^{1 \rightarrow 2} = 0.4$ , and  $p_{1j}^{1 \rightarrow 3} = 0.2 \forall i, j, l$ . The 951  
parameters of the model of the availability of the channels to 952  
the SUs are  $p_j^{NF} = 0.7$  and  $p_j^{FN} = 0.3$ . The length of the time 953  
slot  $\Delta T$  is also  $10^{-2}$  s. Similar parameters are used for the five 954  
SUs to clearly illustrate the performance differences obtained 955  
based on the different strategies. 956

In this simulation, we consider only two scenarios. In sce- 957  
nario 1, all SUs deploy a myopic bidding strategy  $\pi_i^{\text{myopic}}$ ,  $i =$  958  
1, 2, ..., 5, whereas in scenario 2, SU 5 deploys the multiuser 959  
learning-based bidding strategy  $\pi_5^{\mathcal{L}_5}$  with the discount factor 960  
of 0.5, and the other SUs deploy the myopic bidding strategy 961  
 $\pi_i^{\text{myopic}}$ ,  $i = 1, \dots, 4$ . The packet loss rate and cost per time slot 962  
incurred by the SUs are presented in Table III. The accumulated 963  
packet loss and cost of SU 5 for the five scenarios are plotted in 964  
Fig. 6(a) and (b), respectively. 965

Similar to the two-SU network, SU 5 significantly reduces 966  
the packet loss rate by 14.6% and average cost by 16.1% 967  
by adopting the best-response learning-based bidding strategy. 968  
Fig. 6(a) and (b) further verifies the improvement of the per- 969  
formance for SU 1. However, the other SUs' performances are 970  
decreased as they now need to compete against a learning SU 971  
(i.e., SU 5), which is able to make better bids for the available 972  
resources. 973

## B. Multiuser Learning and Delay Impact in a Wireless Test Bed

To validate the performance of multiuser learning and the 976  
impact of various delays in a realistic network setting, we 977  
considered two SUs competing for the available TxOps in our 978  
802.11a-enabled wireless test bed [31]. The channel condition 979  
experienced by the SUs varied between 10 and 30 dB, and 980  
we represented this variation using ten states ( $K = 10$ ). The 981  
parameters of the TxOp model are  $p_j^{NF} = 0.6$  and  $p_j^{FN} = 0.4$ . 982  
The length of the time slot  $\Delta T$  is also  $10^{-2}$  s. The SUs stream 983  
the delay-sensitive video traffic (e.g., the Mobile sequence en- 984  
coded using an H.264 video encoder) to their own destinations 985  
with an average data rate of 1.5 Mb/s. We compare three 986  
scenarios. In scenario 1, both SUs deploy a myopic bidding 987  
strategy  $\pi_i^{\text{myopic}}$ ,  $i = 1, 2$ . In scenario 2, SU 1 deploys the 988  
learning-based bidding strategy  $\pi_1^{\mathcal{L}_1}$  with a discount factor of 989  
0.5, and SU 2 deploys a myopic strategy  $\pi_2^{\text{myopic}}$ . In scenario 3, 990  
both SUs deploy the learning-based bidding strategy  $\pi_i^{\mathcal{L}_i}$ ,  $i =$  991  
1, 2. In the mentioned three scenarios, video applications are 992

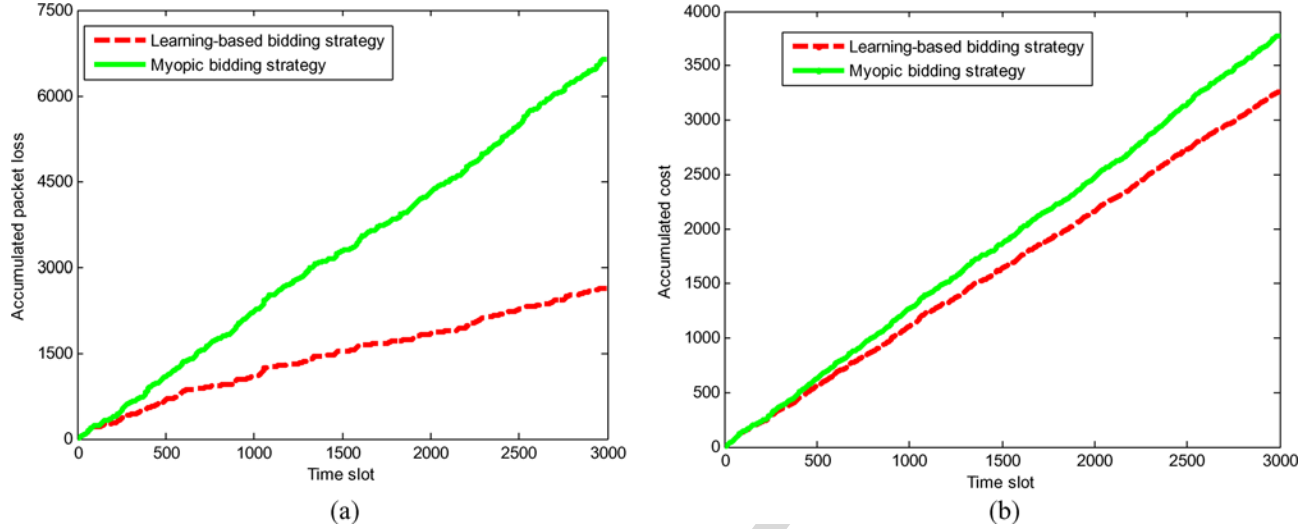


Fig. 6. Accumulated packet loss and cost of SU 5 in the two scenarios. (a) Accumulated packet loss over the time slot. (b) Accumulated cost over the time slot.

TABLE IV  
PERFORMANCE OF SU 1 AND 2 WITH VARIOUS BIDDING STRATEGIES IN THE MORE REALISTIC NETWORK

	Bidding strategies	SU 1		SU 2	
		PSNR (dB)	Average cost	PSNR (dB)	Average cost
Scenario 1	$\pi_1^{\text{myopic}}, \pi_2^{\text{myopic}}$	30.8	5.8951	30.7	5.8845
Scenario 2	$\pi_1^{\mathcal{L}_1}, \pi_2^{\text{myopic}}$	33.0	5.3449	29.9	6.2236
Scenario 3	$\pi_1^{\mathcal{L}_1}, \pi_2^{\mathcal{L}_2}$	31.8	5.6493	31.9	5.6536
Scenario 4	$\pi_1^{\mathcal{L}_1}, \pi_2^{\text{myopic}}$	31.2	7.5439	29.2	6.6748

993 considered to tolerate a delay<sup>8</sup> of 533 ms, which is used in some 994 real-time video streaming applications. In scenario 4, SU 1 995 deploys the learning-based bidding strategy  $\pi_1^{\mathcal{L}_1}$  with a discount 996 factor of 0.5, and SU 2 deploys a myopic strategy  $\pi_2^{\text{myopic}}$ . 997 However, in this scenario, SU 1 streams a video sequence that 998 can only tolerate a delay of 266 ms, which is typical for video 999 conferencing applications.

1000 Table IV shows the average video quality in terms of peak 1001 SNR (PSNR)<sup>9</sup> and incurred cost for both SUs under various 1002 scenarios. Comparing scenario 2 with scenario 1, we observe 1003 that the SU using the learning-based bidding strategy improves 1004 the received video quality by 2.2 dB and reduces the incurred 1005 cost by 9.3%. However, as the performance of SU 1 improves, 1006 this also results in worse performance for SU 2. This observa- 1007 tion is similar to the results in Section VII-A1 and has the same 1008 explanation.

1009 In scenario 3, both SUs deploy the learning-based bidding 1010 strategies and are able to better predict the impact of their 1011 current bidding actions on the future cost based on their ob- 1012 servations. Thus, compared with scenario 1, the performance of 1013 both SUs has improved: SU 1 (SU 2) increases by 1 dB (1.2 dB) 1014 in terms of PSNR and reduces its cost by 4.3% (4.0%). Com- 1015 pared to scenario 2, if SU 2 also deploys the learning-based 1016 approach, then SU 2 also observes its estimated future reward 1017 and will increase its bid, thereby reducing the performance

of SU 1. From Table IV, we note that the PSNR of SU 1 is 1018 decreased by 1.2 dB, whereas the PSNR of SU 2 is increased 1019 by 2 dB. We also observe that the cost of SU 1 is increased by 1020 around 5.6%, whereas the cost of SU 2 is decreased by 9.1%. 1021

1022 In scenario 4, since SU 1 streams a video application with a 1023 lower delay deadline, it has to bid more to ensure that packets 1024 with stringent delay deadline are transmitted to the destination, 1025 and hence, SU 1 incurs a higher transmission cost (41% 1026 increased) compared with scenario 2. Although SU 1 bids 1027 more for the limited available resources, the video quality of 1028 SU 1 is reduced by 1.8 dB due to its stringent delay deadline. 1029 Interestingly, the stringent delay deadline of the SU 1's 1030 application also increases the transmission cost of SU 2 and also 1031 reduces its video quality. This is because the higher bid of SU 1 1032 on limited resources automatically increases the bid of SU 2. 1033

### C. Learning With Imperfect Information

1034 In this section, we consider that SU 1 deploys the learning- 1035 based bidding strategy and SU 2 deploys the myopic strategy. 1036 The environment parameters are the same as in Section VII-B. 1037 To quantify the impact of imperfect information about the 1038 environment on SUs' performance, we assume that SU 1 has the 1039 transition probability of TxOps ( $p_j^{NF} = 0.55$  and  $p_j^{FN} = 0.45$ ), 1040 which is slightly different from the true one (i.e.,  $p_j^{NF} = 0.6$  1041 and  $p_j^{FN} = 0.4$ ). Table V shows the PSNRs and corresponding 1042 cost of both SUs when SU 1 has perfect or imperfect informa- 1043 tion about the TxOps.

1044 From Table V, we observe that an inaccurate model of TxOps 1045 reduces the performance of SU 1 (i.e., the PSNR decreases by

<sup>8</sup>During the simulations, for simplicity, we assume that the packets within one Group of Picture (GOP) have the same delay deadline.

<sup>9</sup>PSNR is a widely adopted metric to objectively measure the video quality. A PSNR difference of 1 dB is significant and can be seen by an untrained human observer.

TABLE V  
PERFORMANCE COMPARISON BETWEEN THE SCENARIOS WHETHER SU 1 HAS PERFECT INFORMATION OR NOT

	Bidding strategies	SU 1		SU 2	
		PSNR (dB)	Average cost	PSNR (dB)	Average cost
Scenario 1 (SU 1 has perfect information)	$\pi_1^{\mathcal{L}}, \pi_2^{myopic}$	33.0	5.3449	30.7	6.2236
Scenario 2 (SU 1 has imperfect information)	$\pi_1^{\mathcal{L}}, \pi_2^{myopic}$	32.7	5.5685	30.5	6.4385

TABLE VI  
CHANNEL AVAILABILITY PROBABILITY

	Channel 1			Channel 2		
	$p_1^{NF}$	$p_1^{FN}$	Number of opportunities	$p_2^{NF}$	$p_2^{FN}$	Number of opportunities
Scenario 1	0.8	0.2	3502	0.8	0.2	3498
Scenario 2	0.5	0.5	2490	0.5	0.5	2462
Scenario 3	0.4	0.6	1960	0.4	0.6	1968

TABLE VII  
AVERAGE PACKET LOSS RATE AND COST FOR THE SUs UNDER VARIOUS RESOURCE CONSTRAINTS

		SU 1		SU 2	
		Packet loss rate	Average cost	Packet loss rate	Average cost
Scenario 1	$\pi_1^{myopic}, \pi_2^{myopic}$	3.08	0.2678	2.90	0.2844
	$\pi_1^{\mathcal{L}}, \pi_2^{myopic}$	2.69	0.3092	4.17	0.4110
Scenario 2	$\pi_1^{myopic}, \pi_2^{myopic}$	21.36	1.8954	23.85	1.7471
	$\pi_1^{\mathcal{L}}, \pi_2^{myopic}$	14.54	1.6764	30.67	2.1744
Scenario 3	$\pi_1^{myopic}, \pi_2^{myopic}$	45.01	3.6283	45.42	3.8289
	$\pi_1^{\mathcal{L}}, \pi_2^{myopic}$	35.21	3.2590	56.44	4.5162

1046 0.3 dB and increases the cost by 4.2%). We further note that this  
1047 will also affect the performance of SU 2. In this simulation, the  
1048 PSNR of SU 2 is reduced by 0.2 dB, and the cost is increased  
1049 by 3.5%. This performance loss can be explained as follows.  
1050 Since SU 1 has an inaccurate model about the available TxOps,  
1051 it may generate a suboptimal bid vector at each stage, which  
1052 will accordingly result in a suboptimal allocation (TxOps and  
1053 payment) among the SUs. This suboptimal allocation will also  
1054 lead to the performance loss of other SUs. Hence, it is essential  
1055 for the users to learn and accurately predict their environment.

#### 1056 D. Impact of Various Dynamics on Learning

1057 In Section VII-A, we demonstrated that the best-response  
1058 learning algorithm improves the bidding strategy, thereby lead-  
1059 ing to a reduced packet loss rate and average cost. In this  
1060 simulation, we further investigate how various dynamics impact  
1061 the learning algorithm proposed in Section VI-D. Specifically,  
1062 we compare the learning performance under different channel  
1063 dynamics, i.e., various available spectrum opportunities for the  
1064 SUs, as discussed in Section II. The source characteristics and  
1065 channel conditions experienced by the SUs are kept the same as  
1066 in Section VII-A1. We consider three types of channel dynam-  
1067 ics corresponding to scenarios 1–3. The transition probabilities  
1068 of TxOps for all three scenarios are listed in Table VI. In each  
1069 scenario, we compare two cases. In the first one, both SUs  
1070 deploy myopic bidding strategies, and in the second one, SU  
1071 1 deploys the best-response learning-based bidding strategy,  
1072 while SU 2 still uses the myopic bidding strategy.

1073 Table VII shows the average packet loss rate and cost ex-  
1074 perenced by the SUs under various channel dynamics. In-

1075 terestingly, we observe from these results that even though  
1076 the learning algorithm reduces the packet loss rate, it does  
1077 not reduce the cost associated with SU 1 when the channel  
1078 resources are abundant as in scenario 1. As the resources  
1079 become increasingly scarce, the learning algorithm helps SU 1  
1080 to simultaneously reduce the packet loss rate and cost, e.g.,  
1081 in scenarios 2 and 3. This observation can be explained as  
1082 follows. When the resources are abundant, the cost (including  
1083 the packet loss and tax) is small, i.e., the “value” of the chan-  
1084 nel is limited, and hence, the learning-based bidding strategy  
1085 does not significantly benefit. On the other hand, when the  
1086 resources are scarce, the bid vectors of the SUs in the current  
1087 time slot will significantly affect the transition of their states  
1088 through the channel allocation compared with the case when  
1089 the resources are abundant. For example, if an SU makes low  
1090 bids as compared to other SUs, it might have no resources  
1091 (channels) allocated to it when resources are scarce (i.e., the SN  
1092 is congested). In this case, the learning-based bidding strategy  
1093 will carefully plan the bid by considering the future impact, and  
1094 thus, it is able to successfully improve the performance of SU 1  
1095 in terms of reducing the average cost.

## 1096 VIII. CONCLUSION AND FUTURE RESEARCH

1097 In this paper, we have modeled the dynamic resource allo-  
1098 cation problem as a “stochastic game” played among strategic  
1099 SUs. At each stage of the game, the CSM deploys a general-  
1100 ized second-price auction mechanism to allocate the available  
1101 spectrum resource. The SUs are allowed to simultaneously  
1102 and independently make bid decisions on that resource by

1103 considering their current states, experienced environment, and  
 1104 estimated future reward. To improve the bid decision at each  
 1105 stage, we propose a best-response learning algorithm to predict  
 1106 the possible future reward at each state. The simulation results  
 1107 show that our proposed learning algorithm can significantly  
 1108 improve the SUs' performance.

1109 We note that the constraint of the perfect information about  
 1110 the available wireless resources can be relaxed for the case  
 1111 when the CSM and wireless users do not have perfect infor-  
 1112 mation about the available resources. In this case, the wire-  
 1113 less users can estimate and build a belief about the available  
 1114 resource. Hence, the stochastic game model can be extended  
 1115 to partially observably stochastic games [32]. This is one of  
 1116 our interesting future research topics. We also note that we  
 1117 can allow the wireless users to adapt their transmission power,  
 1118 which will lead to different interference levels to other users.  
 1119 In this case, the wireless users compete with each other for  
 1120 lower interference levels incurred by other users [6] instead  
 1121 of competing for the transmission time. This can also be for-  
 1122 mulated as a stochastic game, and similar learning algorithms  
 1123 can be developed. This forms another interesting topic of our  
 1124 future research. Our future work also includes analyzing the  
 1125 performance of SNs, where multiple SUs are deploying various  
 1126 learning strategies and protocols.

## 1127 REFERENCES

- 1128 [1] Fed. Commun. Comm., *Spectrum Policy Task Force*, Nov. 2002.  
 1129 Rep. ET Docket No. 02-135.
- 1130 [2] S. Haykin, "Cognitive radio: Brain-empowered wireless communica-  
 1131 tions," *IEEE J. Sel. Areas Commun.*, vol. 23, no. 2, pp. 201–220,  
 1132 Feb. 2005.
- 1133 [3] I. F. Akyildiz, W. Y. Lee, M. C. Vuran, and S. Mohanty, "NeXt  
 1134 generation/dynamic spectrum access/cognitive radio wireless networks:  
 1135 A survey," *Comput. Netw.*, vol. 50, no. 13, pp. 2127–2159, Sep. 2006.
- 1136 [4] C. Kloeck, H. Jaekel, and F. Jondral, "Auction sequence as a new  
 1137 resource allocation mechanism," in *Proc. VTC*, Dallas, TX, Sep. 2005,  
 1138 pp. 240–244.
- 1139 [5] F. Fu, A. R. Fattahi, and M. van der Schaar, "Game-theoretic paradigm  
 1140 for resource management in spectrum agile wireless networks," in *Proc.*  
 1141 *IEEE ICME*, 2006, pp. 873–876.
- 1142 [6] J. Huang, R. Berry, and M. L. Honig, "Auction-based spectrum sharing,"  
 1143 *ACM Mobile Netw. Appl. J. (MONET)*, vol. 11, no. 3, pp. 405–418,  
 1144 Jun. 2006.
- 1145 [7] Y. Xing, R. Chandramouli, and C. M. Cordeiro, "Price dynamics in com-  
 1146 petitive agile spectrum access markets," *IEEE J. Sel. Areas Commun.*,  
 1147 vol. 25, no. 3, pp. 613–621, Apr. 2007.
- 1148 [8] L. Berlemann, S. Mangold, G. R. Hiertz, and B. H. Walke, "Policy defined  
 1149 spectrum sharing and medium access for cognitive radios," *J. Commun.*,  
 1150 vol. 1, no. 1, pp. 1–12, Apr. 2006.
- 1151 [9] C. T. Chou, S. Shankar N, H. Kim, and K. Shin, "What and how much  
 1152 to gain by spectrum agility?" *IEEE J. Sel. Areas Commun.*, vol. 25, no. 3,  
 1153 pp. 576–588, Apr. 2007.
- 1154 [10] S. Shankar, C. T. Chou, K. Challapali, and S. Mangold, "Spectrum agile  
 1155 radio: Capacity and QoS implications of dynamic spectrum assignment,"  
 1156 in *Proc. Global Telecommun. Conf.*, Nov. 2005, pp. 2510–2516.
- 1157 [11] D. Bertsekas and R. Gallager, *Data Networks*. Upper Saddle River, NJ:  
 1158 Prentice-Hall, 1987.
- 1159 [12] M. van der Schaar and S. Shankar, "Cross-layer wireless multimedia  
 1160 transmission: Challenges, principles, and new paradigms," *IEEE Wireless*  
 1161 *Commun.*, vol. 12, no. 4, pp. 50–58, Aug. 2005.
- 1162 [13] *Wireless Medium Access Control (MAC) and Physical Layer (PHY) Spec-*  
 1163 *ifications: Medium Access Control (MAC) Enhancements for Quality of*  
 1164 *Service (QoS), Draft Supplement*, IEEE Std. 802.11e/D5.0, Jun. 2003.
- 1165 [14] R. W. Lucky, "Tragedy of the commons," *IEEE Spectr.*, vol. 43, no. 1,  
 1166 p. 88, Jan. 2006.
- 1167 [15] R. G. Gallager, *Discrete Stochastic Processes*. Norwell, MA: Kluwer,  
 1168 1996.
- [16] L. S. Shapley, "Stochastic games," *Proc. Nat. Acad. Sci. U.S.A.*, vol. 39, 1169  
 no. 10, pp. 1095–1100, Oct. 1953. 1170
- [17] C. Watkins and P. Dayan, "Q-learning, technical note," *Mach. Learn.*, 1171  
 vol. 8, no. 3/4, pp. 279–292, May 1992. 1172
- [18] M. Bowling and M. Veloso, "Rational and convergent learning in stochas- 1173  
 tic games," in *Proc. 17th IJCAI*, Aug. 2001, pp. 1021–1026. 1174
- [19] P. Klemperer, "Auction theory: A guide to the literature," *J. Econ. Surv.*, 1175  
 vol. 13, no. 3, pp. 227–286, Jul. 1999. 1176
- [20] S. P. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, 1177  
 U.K.: Cambridge Univ. Press, 2004. 1178
- [21] F. Fu and M. van der Schaar, "Noncollaborative resource management for 1179  
 wireless multimedia applications using mechanism design," *IEEE Trans.* 1180  
*Multimedia*, vol. 9, no. 4, pp. 851–868, Jun. 2007. 1181
- [22] D. Fudenberg and D. K. Levine, *The Theory of Learning in Games*. 1182  
 Cambridge, MA: MIT Press, 1999. 1183
- [23] M. Jackson, "Mechanism theory," in *Encyclopedia of Life Support* 1184  
*Systems*. Oxford, U.K.: EOLSS, 2003. 1185
- [24] Q. Zhang and S. A. Kassam, "Finite-state Markov model for Rayleigh 1186  
 fading channels," *IEEE Trans. Commun.*, vol. 47, no. 11, pp. 1688–1692, 1187  
 Nov. 1999. 1188
- [25] A. Ortega, "Variable bit-rate video coding," in *Compressed Video Over* 1189  
*Networks*, M.-T. Sun and A. R. Reibman, Eds. New York: Marcel 1190  
 Dekker, 2000, pp. 343–382. 1191
- [26] S. Lal and E. S. Sousa, "Distributed resource allocation for DS-CDMA- 1192  
 based multimedia ad hoc wireless LANs," *IEEE J. Sel. Areas Commun.*, 1193  
 vol. 17, no. 5, pp. 947–967, May 1999. 1194
- [27] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. 1195  
 Cambridge, MA: MIT Press, 1998. 1196
- [28] M. Chiang, S. H. Low, A. R. Calderbank, and J. C. Doyle, "Layering as 1197  
 optimization decomposition: A mathematical theory of network architec- 1198  
 tures," *Proc. IEEE*, vol. 95, no. 1, pp. 255–312, Jan. 2007. 1199
- [29] F. Kelly, A. Maulloo, and D. Tan, "Rate control for communication net- 1200  
 works: Shadow prices, proportional fairness and stability," *J. Oper. Res.* 1201  
*Soc.*, vol. 49, no. 3, pp. 237–252, Mar. 1998. 1202
- [30] X. Zhu, P. Agrawal, J. P. Singh, T. Alpcan, and B. Girod, "Rate allocation 1203  
 for multi-user video streaming over heterogeneous access networks," in 1204  
*Proc. ACM MM*, Sep. 2007, pp. 37–46. 1205
- [31] D. Krishnaswamy and J. Vicente, "Scalable adaptive wireless networks 1206  
 for multimedia in the proactive enterprise," *Intel Technol. J.*, vol. 8, no. 4, 1207  
 pp. 291–301, Nov. 2004. [Online]. Available: [http://developer.intel.com/](http://developer.intel.com/technology/itj/2004/volume08issue04/art04_scalingwireless/p01_abstract) 1208  
[http://developer.intel.com/](http://developer.intel.com/technology/itj/2004/volume08issue04/art04_scalingwireless/p01_abstract) 1209  
[http://developer.intel.com/](http://developer.intel.com/technology/itj/2004/volume08issue04/art04_scalingwireless/p01_abstract) 1210  
[http://developer.intel.com/](http://developer.intel.com/technology/itj/2004/volume08issue04/art04_scalingwireless/p01_abstract) 1211
- [32] D. S. Bernstein, E. A. Hansen, S. Zilberstein, and C. Amato, "Dynamic 1211  
 programming for partially observable stochastic games," in *Proc. AAAI* 1212  
*Spring Symp. Bridging Multi-Agent Multi-Robot. Res. Gap*, Stanford, CA, 1213  
 Mar. 2004. 1214

**Fangwen Fu** received the B.S. and M.S. degrees from Tsinghua University, 1215 AQ3  
 Beijing, China, in 2002 and 2005, respectively. He is currently working toward 1216  
 the Ph.D. degree with the Department of Electrical Engineering, University of 1217  
 California at Los Angeles. 1218

During the summer of 2006, he was an Intern with the IBM T. J. Watson 1219  
 Research Center, Yorktown Heights, NY. His research interests include wireless 1220  
 multimedia streaming, resource management for networks and systems, applied 1221  
 game theory, and video processing and analysis. 1222

**Mihaela van der Schaar** (SM'04) received the Ph.D. degree from the 1223 AQ4  
 Eindhoven University of Technology, Eindhoven, The Netherlands, in 2001. 1224

She is currently an Associate Professor with the Department of Electrical 1225  
 Engineering, University of California, Los Angeles. She is the holder of 1226  
 30 granted U.S. patents and three ISO awards. Her research interests are in 1227  
 multimedia communications, networking, processing, and systems. 1228

Dr. van der Schaar received the National Science Foundation Career Award 1229  
 in 2004, the Best Paper Award from IEEE TRANSACTIONS ON CIRCUITS AND 1230  
 SYSTEMS FOR VIDEO TECHNOLOGY in 2005, the Okawa Foundation Award 1231  
 in 2006, the IBM Faculty Award in 2005, 2007, and 2008, and the Most Cited 1232  
 Paper Award from *EURASIP: Image Communications* in 2006. 1233