

# Wireless Communications and Mobile Computing

## Finder-MCTS: A Cognitive Spectrum Allocation Based on Traveling State Priority and Scenario Simulation in IoV

Zhong Li<sup>1</sup>, Hao Shao<sup>1</sup>

<sup>1</sup>College of Information Science and Technology, Donghua University, Shanghai 201620, China

Correspondence should be addressed to Hao Shao; shao0216@mail.dhu.edu.cn

### Abstract

With the increasing number of intelligent connected vehicles, the problem of scarcity of communication resources has become increasingly obvious. It is a practical issue with important significance to explore a real-time and reliable dynamic spectrum allocation scheme for the vehicle users, while improving the utilization of available spectrum. However, previous studies have problems such as local optimum, complex parameter setting, learning speed, and poor convergence. Thus, in this paper, we propose a cognitive spectrum allocation method based on traveling state priority and scenario simulation in IoV, named Finder-MCTS. The proposed method integrates offline learning with online search. This method mainly consists of two stages. Initially, Finder-MCTS gives the allocation priority of different vehicle users based on the vehicle's local driving status and global communication status. Furthermore, Finder-MCTS can search for the approximate optimal allocation solutions quickly online according to the priority and the scenario simulation, while with the offline deep neural network based environmental state predictor. In the experiment, we use SUMO to simulate the real traffic flows. Numerical results show that our proposed Finder-MCTS has 36.47%, 18.24%, 9.00% improvement on average than other popular methods in convergence time, link capacity and channel utilization, respectively. In addition, we verified the effectiveness and advantages of Finder-MCTS compared with two MCTS algorithms' variations.

**Keywords:** Internet of Vehicle (IoV); cognitive radio; dynamic spectrum allocation; Monte-Carlo tree search (MCTS)

## 1. Introduction

Recently, as a promising technology, internet of vehicles (IoV) has attracted the attention of governments and enterprises around the world, to serve the smart city. The moving vehicles can be regarded as mobile terminals equipped with advanced network components, such as wireless network interfaces, on-board sensors, which provide many personalized services by accessing the internet. These vehicle services (*e.g.*, road condition broadcasts, dangerous event predictions) have high requirements for data transmission and communication quality. Although 5G technology is becoming popular and growing rapidly, the available spectrum resources have not increased simultaneously. So far, the spectrum resources of 6GHz and below 6GHz have almost been exhausted [1]. Moreover, the spectrum resources at the base stations are usually allocated to the calls and traffic services of mobile phone users first. Thus, the scarcity of spectrum resources and the low utilization of frequency bands are critical issues hindering the development of IoV.

33 Currently, as an effective solution to the underutilized problem of spectrum resources, cognitive radio (CR)  
 34 can reuse the idle spectrum resources through dynamic spectrum access technology. In CR networks, network  
 35 users are divided into two types: primary users (PUs) and secondary users (SUs). PUs have the high priority to  
 36 use the spectrum in the authorized frequency bands. SUs can dynamically access spectrum holes opportunistically  
 37 and use available spectrum resources, which can enhance the spectrum utilization. Therefore, inspired by the  
 38 CR technology, we let vehicles equipped with CR functions and form a cognitive radio-based internet of vehicles  
 39 (CR-IoV). We utilize CR to help solve the low utilization of frequency bands in IoV.

40 In CR-IoV, the system includes PUs (composed by mobile phone users) and SUs (composed by vehicles  
 41 equipped with CR functions). However, in reality, vehicle users with high mobility will cause frequent changes  
 42 in the network topology. The availability of spectrum will also change with the activation time and channel oc-  
 43 cupancy of PUs. Hence, how to meet the real-time and reliable requirements when solving the dynamic spectrum  
 44 allocation problem under a time-varying environment is an significant challenge.

45 There are many previous studies about dynamic spectrum allocation in mobile wireless networks. The most  
 46 popular studies can be mainly classified into four categories: (1) traditional optimization theory-based allocation  
 47 methods [2, 3]; (2) game theory-based allocation methods [4–6]; (3) swarm intelligence optimization-based allo-  
 48 cation methods [7–11]; (4) machine learning-based allocation methods [12–17]. Although the above methods can  
 49 solve the spectrum allocation problem, there exist many disadvantages. First, when the constraints are complex,  
 50 traditional optimization theory and game theory are not suitable for quickly solving the large-scale dynamic plan-  
 51 ning problems. Second, the swarm intelligence optimization is easy to fall into the local optimum [18]. Besides, the  
 52 effective parameter settings and selection in the swarm intelligence optimization is also complex. Recently, deep  
 53 reinforcement learning (DRL) algorithms have been proved to solve complex dynamic decision-making problem  
 54 with high-dimensional state and action space. It can learn the potential regularities in the environment with the  
 55 help of the idea of trial and error, thereby assisting the intelligent decision-making. However, this type of machine  
 56 learning-based method also exists some limitations, such as slow learning speed, poor convergence, and bad self-  
 57 adaptation ability. Thus, in this paper, we propose a new cognitive spectrum allocation method based on traveling  
 58 state priority and different scenarios specially for IoV in this paper.

59 First, especially in IoV, we should consider the traveling/moving state of a vehicle. A vehicle that is about to  
 60 leave the coverage area of a base station should have relatively low spectrum allocation priority. Vehicle users with  
 61 different traveling state, such as location, speed, acceleration, communication capabilities, should have different  
 62 opportunities to obtain spectrum resources. Thus, in this paper, we consider the priority assignment based on  
 63 vehicle traveling state when doing spectrum allocation.

64 In addition, in this proposed new method, we choose Monte-Carlo tree search algorithm (MCTS) to model  
 65 our problem. Traditional model-free based deep reinforcement learning algorithms (*e.g.*, deep Q network, soft  
 66 actor-critic) often require a large amount of samplings and learn strategies from past experiences with the help of  
 67 neural networks. However, model-based deep MCTS can not only use deep neural networks to fit the environment  
 68 model from experience data, but also can simulate a variety of possible future trajectories for evaluation through the  
 69 expansion of the tree structure, so as to choose more promising directions to explore the best policy. In this paper,  
 70 through designing to simulate different scenarios, we improve the learning efficiency and reduce the searching  
 71 space compared with traditional MCTS methods.

72 Our main contributions can be summarized as follows:

73 • We design a priority assignment rule based on vehicle traveling state for spectrum allocation. Through defin-  
 74 ing a vehicle traveling evaluation score and a network utility score, we obtain a comprehensive priority evaluation  
 75 score for each vehicle. According to the priority score, we allocate available spectrum resources from the highest  
 76 priority to the lowest vehicle user, which can improve the allocation performance when doing dynamic spectrum  
 77 allocation in IoV.

78 • Combining with the above priority score, we propose a cognitive spectrum allocation method based on travel-  
 79 ing state priority and different scenarios specially for IoV, named Finder-MCTS. We model the problem of spectrum  
 80 allocation as a binary integer linear programming problem (BILP) with constraints. Meanwhile, through designing

81 a constraint oriented tree expansion and scenario simulation mechanism, Finder-MCTS can give an approximate  
 82 optimal solution quickly and improve the link capacity of V2I (vehicle to infrastructure) communication in the  
 83 network.

84 • We conduct experiments to evaluate the performance of Finder-MCTS by using SUMO. Results show that our  
 85 proposed method has 36.47%, 18.24%, 9.00% improvements on average than other popular comparison methods in  
 86 convergence time, link capacity and channel utilization, respectively. In addition, Finder-MCTS also shows good  
 87 improvements with the aid of priority evaluation and different scenarios' simulation of PUs' service durations,  
 88 compared with two variations of MCTS.

89 The remainder of this paper is organized as follows. In Section 2, a review of related work is provided. In Sec-  
 90 tion 3, the system scenario and problem formalization are presented in detail. In Section 4, the priority assignment  
 91 based on vehicle traveling state are described. In Section 5, the Finder-MCTS method for cognitive IoV spectrum  
 92 allocation are proposed. In Section 6, simulations are carried out to demonstrate the effectiveness of the proposed  
 93 Finder-MCTS method. In Section 7, conclusion and future work are given.

## 94 2. Related Work

95 Nowadays, there are many excellent studies on dynamic spectrum allocation in cognitive radio networks. In this  
 96 section, we classify and compare them from the perspective of theoretical methods.

### 97 2.1 Spectrum Resource Allocation Based on Traditional Optimization Theory and Game 98 Theory

99 In order to solve the problem of dynamic allocation of spectrum resources in wireless communications, the tradi-  
 100 tional methods mainly include the methods based on mathematical optimization [2, 3] and the methods based on  
 101 game theory [4–6]. For example, Martinovic *et al.* propose a cognitive radio spectrum allocation method based on  
 102 integer linear programming in the work of [3], which solves the spectrum allocation problem with interference by  
 103 using many complex assumptions and constraints. It is difficult or even impossible to find an optimal solution in the  
 104 real cognitive radio network with the complex environment and dynamic network topology. Although the methods  
 105 based on mathematical optimization have high solution accuracy, the generalization capability is insufficient.

106 Besides, with the goal of maximizing spectrum utilization, Yi *et al.* introduce a spectrum resource allocation  
 107 method based on auction in the work of [5]. Liu *et al.* design a dynamic spectrum access method using game theory  
 108 in the work of [6]. However, these methods are not fit for IoV. The high mobility of vehicles puts forward a strict  
 109 requirement for the convergence of Nash equilibrium in the game theory. It is hard to reach this equilibrium point.

### 110 2.2 Spectrum Resource Allocation Based on Swarm Intelligence Optimization

111 There are many related studies [7–11] based on swarm intelligence optimization in the domain of spectrum allo-  
 112 cation. For example, Liu *et al.* use PSO to solve the allocation of spectrum resources in a centralized way in the  
 113 work of [11]. However, the iteration of swarm intelligence optimization usually gets stuck in local optimal solu-  
 114 tions, which can be far from the global optimal solution [18]. In addition, many swarm intelligence optimization  
 115 algorithms have large amount of calculations in debugging due to complex parameters.

### 116 2.3 Spectrum Resource Allocation Based on Machine Learning

117 In recent years, with the development of statistical learning methods, many studies use machine learning to realize  
 118 the dynamic spectrum allocation [12]. Among them, reinforcement learning can guide a system agent to learn the  
 119 unknown environment by trial and error [19]. It can be applied to the spectrum allocation decision.

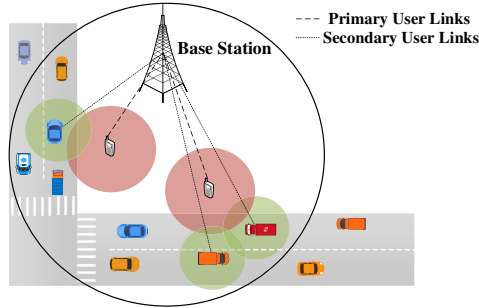


Figure 1: System scenario of spectrum allocation in CR-IoV.

120 First, the multi-arm gambling machine (MAB) is not only an important random decision-making theory in the  
 121 field of operational research, but also belongs to a type of online learning algorithm in reinforcement learning. The  
 122 task of the agent is to select one arm to pull in each round based on the historical rewards it collected, and the goal  
 123 is to collect cumulative reward over multiple rounds as much as possible. In essence, MAB is a way to optimize  
 124 the reward by balancing exploration and exploitation. Li *et al.* give a survey of spectrum resource allocation by  
 125 using MAB in the cognitive radio network in the work of [13]. Zhang *et al.* formulate and study a multi-user MAB  
 126 problem that exploits the idea of temporal-spatial spectrum reuse in the cognitive radio network [14]. However,  
 127 the MAB modeling does not consider the cost of pulling arms in the existing allocation schemes. When MAB is  
 128 utilized to solve the allocation problem in a centralized way, the scale of the arm increases exponentially with the  
 129 number of users to be assigned. Therefore, the convergence of spectrum resource scheduling algorithm based on  
 130 MAB cannot be guaranteed.

131 In addition, model-free-based deep reinforcement learning is also applied to the research of spectrum allocation.  
 132 Naparstek *et al.* propose a spectrum allocation scheme based on deep learning framework under the wireless  
 133 environment in the work of [15]. However, model-free-based deep reinforcement learning has problems of slow  
 134 online learning speed and bad self-adaption ability.

135 Recently, another kind of model-based reinforcement learning, Monte-Carlo tree search algorithm (MCTS), is  
 136 applied in the field of resource allocation [16, 17]. The MCTS-based allocation algorithm builds a decision tree to  
 137 explore the possible solutions by expanding and pruning. Due to the expansion of the tree, the search space becomes  
 138 tremendous gradually and the calculation scale is unacceptable. If this type of method is applied in IoV directly,  
 139 the dynamic environment will further cause a large search tree. In addition, due to the neglect of environmental  
 140 uncertainties, the random strategy adopted by Basic-MCTS in the simulation stage will produce a high variance,  
 141 which reduces the search efficiency [20].

### 142 3. System Scenario and Problem Formalization

143 In this section, we introduce the system scenario of spectrum allocation in CR-IoV in Section , and give the math-  
 144 ematical formalization of our optimization problem in Section .

#### 145 3.1 System Scenario

146 Figure 1 shows the system scenario of spectrum allocation in CR-IoV. PUs are the authorized mobile phone users  
 147 in the current network, and SUs are the vehicles equipped with CR modules. When a PU occupies a channel, there  
 148 is a protection area around the PU (*i.e.*, the red area in Figure 1). Similarly, an interference radius is also generated  
 149 when the SU occupies a channel (*i.e.*, the green area in Figure 1). Any radiation from SUs falling into the protection  
 150 area would interfere with the PU.

151 In this scenario, our designed allocation algorithm is deployed on the base station. Vehicle nodes equipped with  
 152 CR modules can sense whether there exist available idle spectrum resources. A vehicle can use the common control  
 153 channel (CCC) to send a request to the base station to access the channel. The base station collects requests from  
 154 vehicles centrally, and learns a near optimal policy to allocate available channel resources to cognitive vehicles  
 155 within the coverage area (*i.e.*, the black solid circle in Figure 1).

156 Note that, because IoV is a dynamic network, our designed spectrum resources allocation algorithm must  
 157 be executed within a defined allocation time window. We assume that the allocation time window for channel  
 158 allocation is  $T$ . After the time window slides, we will refresh and observe the current vehicles which require to  
 159 access the base station. A large time window cannot meet the real-time requirement of IoV, but a time window that  
 160 is too small cannot support our algorithm for well operating. In the experiment, we set the size of a time window  $T$   
 161 to 10s to handle the dynamic network.

## 162 3.2 Definitions and Problem Formalization

### 163 1) Definitions

164 In this paper, we consider spectrum resource allocation in the underlay mode, *i.e.*, each channel can support the  
 165 parallel transmissions of several access users. Assume that within a base station's communication coverage, there  
 166 are  $N$  SUs competing for spectrum resources of  $M$  channels at time  $t$ , and the channels are orthogonal and non-  
 167 overlapping. Meanwhile, we assume that there are  $K$  PUs as a prerequisite for spectrum allocation in the coverage  
 168 area, and each PU occupies only one channel for information transmission in the current network. The spectrum  
 169 resource allocation model consists of a channel availability matrix  $L$ , a SU-SU interference constraint matrix  $C$ , a  
 170 channel reward matrix  $R$  and a conflict-free channel assignment matrix  $A$ .

171 We define that a PU  $k$  ( $1 \leq k \leq K$ ) occupying a certain authorized channel  $m$  ( $1 \leq m \leq M$ ) in CR-IoV has  
 172 a protection radius  $\tilde{R}(k, m)$ . Meanwhile, each SU  $n$  ( $1 \leq n \leq N$ ) has an interference radius  $\bar{R}(n, m)$  on channel  
 173  $m$  due to its transmit power. We obtain an Euclidean distance  $\mathbb{D}(\tilde{R}(k, m), \bar{R}(n, m))$  between a PU  $k$  and a SU  $n$ .  
 174 When the inequality  $\mathbb{D}(\tilde{R}(k, m), \bar{R}(n, m)) - \tilde{R}(k, m) \leq \bar{R}(n, m)$  holds, it means that there exists communication  
 175 interference between the PU  $k$  and SU  $n$ .

176 Similarly, we also can obtain an Euclidean distance  $\mathbb{D}(\bar{R}(n, m), \bar{R}(n', m))$  between two different SUs  $n$   
 177 and  $n'$ , where  $\bar{R}(n, m)$  and  $\bar{R}(n', m)$  are the interference radius values of the two SUs  $n$  and  $n'$ . When  
 178  $\mathbb{D}(\bar{R}(n, m), \bar{R}(n', m)) - \bar{R}(n, m) \leq \bar{R}(n', m)$  holds, it means that there exists communication interference  
 179 between the two SUs  $n$  and  $n'$ . Note that, when there is no communication interference between two users, they  
 180 can use the same channel for transmissions at the same time; otherwise, they cannot access the same channel at the  
 181 same time.

182 Next, according to the above descriptions of communication interference between different users, we give the  
 183 following definitions about our problem.

#### 184 • Channel Availability Matrix $L$ .

$L = \{l_{n,m} | l_{n,m} \in \{0, 1\}\}_{N \times M}$  is an  $N \times M$  dimensional matrix used to describe the channel availability.  
 When  $l_{n,m} = 1$ , it means channel  $m$  is available for SU  $n$ , and vice versa. It needs to meet the following two  
 conditions to determine whether channel  $m$  is available for SU  $n$ . First, SU  $n$  cannot use channel  $m$  occupied  
 by PU  $k$  under the condition  $\mathbb{D}(\tilde{R}(k, m), \bar{R}(n, m)) - \tilde{R}(k, m) \leq \bar{R}(n, m)$ . Second, SUs need to compare the  
 interference power they received with the maximum allowable interference level  $\gamma_m$  on channel  $m$ . Channel  $m$  is  
 considered to be available to SU  $n$  if the following inequality is satisfied:

$$\sum_{k=1}^K P_{m,n,k} + N_m \leq \gamma_m \quad (1)$$

185 where  $P_{m,n,k}$  denotes the received power at SU  $n$  of a signal transmitted from PU  $k$  on channel  $m$ ;  $N_m$  denotes  
 186 the level of background noise on channel  $m$ .

#### 187 • SU-SU Interference Matrix $C$ .

188  $C = \{c_{n,n',m} | c_{n,n',m} \in \{0, 1\}\}_{N \times N \times M}$  is an  $N \times N \times M$  dimensional matrix used to describe the interference  
 189 constraint between two different SUs  $n$  and  $n'$  on channel  $m$ , where  $c_{n,n',m} = 1$  indicates that there exists inter-  
 190 ference when SUs  $n$  and  $n'$  share the channel  $m$  for information transmission. Conversely,  $c_{n,n',m} = 0$  indicates  
 191 that SUs  $n$  and  $n'$  can use channel  $m$  simultaneously. When  $n = n'$ ,  $c_{n,n',m} = 1 - l_{n,m}$ . Meanwhile, the matrix  
 192 element needs to satisfy the condition  $c_{n,n',m} \leq l_{n,m} \cdot l_{n',m}$ , *i.e.*, the premise for the possibility of interference is  
 193 that channel  $m$  is available to both SUs  $n$  and  $n'$ .

194 • **Channel Allocation Matrix  $A$ .**

195  $A = \{a_{n,m} | a_{n,m} \in \{0, 1\}\}_{N \times M}$  is an  $N \times M$  dimensional matrix used to describe the conflict-free channel  
 196 allocation for SUs. When  $a_{n,m} = 1$  holds, it means that the channel  $m$  is allocated to the SU  $n$ , and vice versa.  
 197 Meanwhile, matrix  $A$  must satisfy the interference constraint given by matrix  $C$ . That is to say, for two different  
 198 SUs  $n$  and  $n'$ , when  $c_{n,n',m} = 1$ , the equation  $a_{n,m} \cdot a_{n',m} = 0$  holds. In addition, we assume that each SU in the  
 199 allocation can only occupy one channel for information transmission. Therefore, for any two different channels  $m$   
 200 and  $m'$ , the inequality  $a_{n,m} + a_{n,m'} \leq 1$  should be satisfied.

201 • **Channel Reward Matrix  $R$ .**

$R = \{r_{n,m} | r_{n,m} \geq 0\}_{N \times M}$  is an  $N \times M$  dimensional matrix used to describe the link rewards for different  
 SUs. Notation  $r_{n,m}$  denote the reward obtained by SU  $n$  when it occupies channel  $m$  of a base station.  $r_{n,m}$  is  
 measured by the link capacity. Link capacity is defined as follows:

$$r_{n,m} = W_m \cdot \log_2(1 + SINR_{n,m}) \quad (2)$$

where  $W_m$  is the bandwidth of channel  $m$ , and  $SINR_{n,m}$  is the signal to interference plus noise ratio when SU  $n$   
 accesses channel  $m$ . The calculation of  $SINR_{n,m}$  is shown in Eq. (3).

$$SINR_{n,m} = \frac{P_{m,n}}{N_m + \sum_{q=1, q \neq n}^{Count(A_m)} P_{m,q}} \quad (3)$$

202 where we regard the SU and the base station as the transmitting-end and the receiving-end respectively. Here,  $A_m$   
 203 represents the  $m$ -th column vectors of matrices  $A$  and  $Count(A_m)$  denotes the total number of allocated SUs on  
 204 channel  $m$ ;  $P_{m,n}$  is the power received by the receiver (base station) from the transmitter  $n$  on channel  $m$ .

205 **2) Problem Formalization**

From the above definitions, it can be seen that there are more than one channel allocation matrices satisfying the  
 allocation constraints. Therefore, let  $\Lambda_{L,C} = \{A_g\} (g \in \mathbb{N}^+)$  denote the set of all conflict-free channel allocation  
 schemes derived from the current network conditions  $L$  and  $C$ . Because there are many possible spectrum allocation  
 schemes, choosing different spectrum allocation schemes will generate different total system rewards. The object  
 of spectrum allocation in this paper is to maximize the total network capacity  $U(A, R)$  of the network system. We  
 give the definition of total network capacity as follows:

$$U(A, R) = SUM\left(\sum_{m=1}^M A_m \circ R_m\right) \quad (4)$$

206 where  $R_m$  represents the  $m$ -th column vector of matrix  $R$ . Notation  $\circ$  represents the Hadamard product, *i.e.*,  
 207 multiplication of the elements at the corresponding positions of the two vectors.  $A_m$  is a 0/1 decision vector of  
 208  $N \times 1$  size, and  $R_m$  is an  $N \times 1$  dimensional reward vector with real numbers.  $\sum_{m=1}^M A_m \circ R_m$  is also an  $N \times 1$   
 209 dimensional vector. Notation  $SUM$  is the operator that returns the summation of all entries of a matrix.

210 In the IoV, our paper aims to obtain an optimal channel allocation matrix  $A^*$  (*i.e.*, with the equation  $A^* =$   
 211  $\underset{A \in \Lambda_{L,C}}{argmax} U(A, R)$ ), which satisfies the above non-interference constraints and solves the problem of low utilization  
 212 of spectrum resources at the base station side. The combinatorial optimization problem can be formulated as a  
 213 binary integer linear programming problem (BILP) as follows:

$$\max_{A,R} U(A, R) = \max_{A,R} SUM\left(\sum_{m=1}^M A_m \circ R_m\right) \quad (5)$$

s.t. :

$$a_{n,m} \in \{0, 1\}, r_{n,m} \geq 0, 1 \leq n \leq N, 1 \leq m \leq M \quad (5a)$$

$$a_{n,m} \leq l_{n,m} \quad (5b)$$

$$a_{n,m} \cdot a_{n',m} = 0 \quad \text{if } c_{n,n',m} = 1 \quad (5c)$$

$$a_{n,m} + a_{n,m'} \leq 1 \quad (5d)$$

$$\bar{P}_{m,n}^{min} \leq \bar{P}_{m,n} \leq \bar{P}_{m,n}^{max} \quad \text{if } \bar{P}_{m,n} \neq 0 \quad (5e)$$

$$\sum_{i=1}^n P_{m,k,n} \leq \delta_{m,k}, 1 \leq k \leq K \quad (5f)$$

$$A_m \times R_m^T \leq \phi_m \quad (5g)$$

214 Among these, constraint (5a) gives the value range of the matrix vectors  $A_m$  and  $R_m$ . Constraint (5b) ensures  
 215 that an allocated channel must be an available channel for SU  $n$ . Besides, to protect the communication of each SU  
 216 from interference by other SUs on channel  $m$ , the conflict-free channel allocation matrix should satisfy constraint  
 217 (5c). Constraint (5d) indicates that each SU can only occupy one channel for information transmission. In the  
 218 constraint (5e),  $\bar{P}_{m,n}$  represents the transmission power of SU  $n$  on channel  $m$ ;  $\bar{P}_{m,n}^{min}$  and  $\bar{P}_{m,n}^{max}$  represent the  
 219 maximum and minimum allowable transmission power of SU  $n$  on channel  $m$  respectively. This constraint defines  
 220 the upper and lower bounds for the transmitted power of the SU. In other words, the transmission power of the SU  
 221 should meet two constraints: on one hand, it should not interfere with the normal use of the PU; on the other hand,  
 222 it should meet the minimum allowable SINR required for transmissions. In the constraint (5f),  $P_{m,k,n}$  represents  
 223 the interference power of SU  $n$  received by PU  $k$  on channel  $m$ , and  $\delta_{m,k}$  represents the maximum allowable  
 224 interference power of PU  $k$  on channel  $m$ . For any PU  $k$ , the total received interference power on the channel  
 225  $m$  must be kept below the maximum allowable interference threshold, *i.e.*, the PU is not interfered by SUs on  
 226 the channel. In the constraint (5g),  $\phi_m$  represents the available bandwidth of channel  $m$ , and  $R_m^T$  represents the  
 227 transposed vector of  $R_m$ . This constraint ensures that the total network capacity of channel  $m$  should be less than  
 228 or equal to its available bandwidth.

## 229 4. Priority Assignment Based on Vehicle Traveling State

230 In Section , we describe the problem of priority assignment. In Section , we give the detailed definition of priority.

### 231 4.1 Problem Description

232 In CR-IoV, when the system carries out the spectrum allocation, the current state of vehicle traveling should be  
 233 considered. For example, if a vehicle is about to leave the communication range of the current base station, it  
 234 should be assigned to a low priority for spectrum allocation.

235 The traveling state of a vehicle at the current moment mainly includes direction, speed, acceleration and GPS  
 236 coordinates. Besides, the traveling state also should considers the degree of geographical dispersion among vehicles  
 237 and the communication capability of a vehicle.

238 The current state information of each vehicle is collected by the current communicating base station. Then we  
 239 carry out priority evaluation for different cognitive vehicle users to distinguish the priority weights for spectrum  
 240 allocation.

241 For a SU  $n$  who initiates a service request, from the perspectives of the global state and local state, a compre-  
 242 hensive priority evaluation score  $Priorityscore_n$  is constructed by defining a vehicle traveling evaluation score  
 243  $Travelingscore_n$  and a network utility score  $Utility_n$  for the SU.

## 244 4.2 Priority Definition Based on Vehicle Traveling State

245 *Definition 1: Vehicle Traveling Evaluation Score*

246 According to the GPS coordinates, speed and acceleration, we define a vehicle traveling evaluation score  
 247  $Travelingscore_n$  for a cognitive vehicle  $n$  as

$$Travelingscore_n = \frac{1 + \cos(\theta_n)}{4} \cdot \left( \frac{v_{max} - v_n}{v_{max} - v_{min}} + \frac{1}{1 + e^{a_n}} \right) \quad (6)$$

248 where  $\theta_n$  denotes the angle between the current driving direction and the link connected the vehicle's position with  
 249 the base station's position. Notation  $a_n$  denotes the acceleration of the vehicle  $n$ . Notation  $v_n$  denotes the speed  
 250 of the vehicle  $n$ . Notations  $v_{max}$  and  $v_{min}$  represent the maximum and minimum values of the diving speed. We  
 251 assume that the vehicle speed is within the value interval  $[v_{min}, v_{max}]$ .

252 Obviously, a relatively large angle  $\theta_n$  indicates that vehicle  $n$  will travel out of the coverage range of the base  
 253 station in the future. Therefore, vehicle  $n$  with large  $\theta_n$  should be given a relatively low spectrum allocation priority.  
 254 We use formula  $\frac{1 + \cos(\theta_n)}{2}$  to normalize the different weights of the angle  $\theta_n$  to the value interval  $[0, 1]$ . In addition,  
 255 a vehicle with high driving speed will quickly travel out of the coverage range of the base station in the future.  
 256 Therefore, it should be given a relatively low spectrum allocation priority. The normalized formula  $\frac{v_{max} - v_n}{v_{max} - v_{min}}$   
 257 is used to describe the influence of vehicle driving speed on the priority. Similarly, a vehicle with high acceleration  
 258 should be given a relatively low spectrum allocation priority. To normalize the value interval to  $[0, 1]$ , formula  
 259  $\frac{1}{1 + e^{a_n}}$  is used to describe the influence of vehicle driving acceleration on the priority. Finally, to constrain the value  
 260 of  $Travelingscore_n$  within the value interval  $[0, 1]$ , we use constant coefficient  $\frac{1}{2}$  to obtain the right side of Eq.  
 261 (6).

262 *Definition 2: Network Utility Score*

263 We define a network utility score to evaluate the communication capability of cognitive vehicles. For cognitive  
 264 vehicle  $n$ , its network utility score is defined as follows:

$$Utility_n = \log_2(1 + SNR_n) \cdot \frac{\sum_{1 \leq n, n' \leq N, n' \neq n} Dispersion_{n, n'}}{N - 1} \quad (7)$$

265 where  $SNR_n$  denotes the signal-to-noise ratio of the user  $n$  to receive the signal from the base station. Formula  
 266  $\sum_{1 \leq n, n' \leq N, n' \neq n} Dispersion_{n, n'}$  represents the global dispersion of user  $n$  within the coverage area of the base  
 267 station.

268 For the numerator of Eq. (7), we give the following detailed definition. The  $Dispersion_{n, n'}$  between two SUs  
 269  $n$  and  $n'$  is defined as follows:

$$Dispersion_{n, n'} = \begin{cases} 1 & D_{n, n'} > \varepsilon_n \\ 0 & \text{others} \end{cases} \quad (8)$$

270 where  $\varepsilon_n$  is a dispersion threshold; notation  $D_{n, n'}$  represents the average dispersion time between two SUs  $n$  and  
 271  $n'$ . First, the threshold  $\varepsilon_n$  is obtained by taking the median value of  $\{D_{n, n'} | 1 \leq n' \leq N, n' \neq n\}$ . Second, the  
 272 average dispersion time  $D_{n, n'}$  is defined as

$$D_{n, n'} = \frac{\int_0^T \beta_{n, n'}(t) dt}{\tau_{n, n'}} \quad (9)$$

273 In Eq. (9), the communication dispersion state between two vehicles  $n$  and  $n'$  is defined as  $\beta_{n, n'}(t)$ . When  
 274 there exists communication interference between vehicle  $n$  and  $n'$ , we let  $\beta_{n, n'}(t) = 0$ . It means that the two are in  
 275 an 'encounter' state. On the contrary, when  $\beta_{n, n'}(t) = 1$ , it means that the two are in a 'scattered' state. Thus, in a



time window  $T$ , the numerator of Eq. (9) represents the total dispersion time between user  $n$  and user  $n'$ . Besides,  $\tau_{n,n'}$  in the denominator denotes the total number of times that user  $n$  and user  $n'$  are in the ‘scattered’ state in time window  $T$ . Obviously, the higher the value of  $D_{n,n'}$ , the longer the time that the two users  $n$  and  $n'$  are in the ‘scattered’ state. Thus, we conclude that the higher the global dispersion  $\sum_{1 \leq n, n' \leq N, n' \neq n} Dispersion_{n,n'}$ , the greater the probability that vehicle  $n$  has the chance to reuse the channel, which further leads to a high network utility.

To sum up, a vehicle with a large network utility score in Eq. (7) means that its global communication capability is strong, so the vehicle should be given a high spectrum allocation priority.

*Definition 3: Comprehensive Priority Evaluation Score*

According to the vehicle traveling evaluation score  $TravelingScore_n$  and network utility score  $Utility_n$ , we construct a comprehensive priority evaluation score  $PriorityScore_n$  for the cognitive vehicle  $n$  below,

$$Priorityscore_n = Travelingscore_n \cdot Utility_n \quad (10)$$

For a cognitive vehicle who requests to access the base station, the base station calculates the priority score by collecting the vehicle’s information. We rank all the scores from the largest to the smallest. Therefore, we can obtain a priority order list  $Priorityscore\_list$  for all the cognitive vehicles in the current allocation task, which will be used in the following Section .

## 5. Finder-MCTS Algorithm for Cognitive IoV Spectrum Allocation

In the introduction, we mentioned that our paper will use MCTS to solve the problem of efficient spectrum allocation for CR-IoV. MCTS is a classic reinforcement learning algorithm based on tree search. To distinguish it from the method proposed in our paper, we call the classic MCTS as Basic-MCTS. The Basic-MCTS offers a concise computation framework by recursively using a tree policy to expand the search tree towards high-reward nodes, and a default policy to perform the simulations for updating the estimated rewards and other statistics [21]. However, due to the continuous expansion of search actions, the search scale of Basic-MCTS is often very large, which greatly affects its search speed. In addition, due to the neglect of environmental uncertainties, the random strategy adopted by Basic-MCTS in the simulation stage will produce a high variance, which reduces the search effect of Basic-MCTS.

To improve the search speed and obtain a near optimal solution, we propose an algorithm named Finder-MCTS in this section. First, we construct a search tree vertically according to the comprehensive priority evaluation score defined in Definition 3 above. Meanwhile, the constraints defined in Section are also considered to reduce the search scale of the tree horizontally. Second, the uncertainty of the SUs’ spectrum occupation activities are included into the simulation strategy. We give the bias estimation of reward in different scenarios in the simulation stage so as to approximate the real environment and accelerate the convergence of tree search.

Thus, in Finder-MCTS, the first step is to use Markov decision process (MDP) to construct Monte-Carlo tree computation framework (Section ). Then, with respect to the state prediction, we give a DNN-based environment state predictor–ESP (Section ). Finally, we describe the detailed steps of Finder-MCTS algorithm (Section ).

### 5.1 Finder-MCTS’s Computation Framework

The problems solved by the MCTS are commonly formalized by the Markov decision process (MDP), in which we take the base station as the spectrum scheduling agent and use the link capacity formulated in Eq. (2) as the value of the reward  $Q$  when a SU occupies a channel. Let  $\mathcal{S}$  and  $\mathcal{A}$  denote the MDP state space and action space, respectively.  $\mathcal{F} : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  denotes the MDP transition function from a state-action pair to the next state. The state transition function  $f_{ESP}$  is given by a deep neural network (DNN) simulator in Section . The definitions of the MDP state space and action space are described as follows,

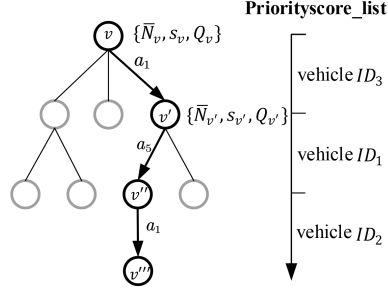


Figure 2: An example of search steps in Finder-MCTS.

$$\mathcal{S} = \{s_v | s_v, \varphi_v, \xi_v\} \quad (11)$$

$$\mathcal{A} = \{a_m | 1 \leq m \leq M\} \quad (12)$$

317

318 In Eq. (11), the MDP state is composed of two parts:  $\lambda_v$  denotes a vector of remaining bandwidth of  $M$   
 319 channels under the base station, with  $s_v = (\lambda_v = (\bar{\lambda}_1, \dots, \bar{\lambda}_m, \dots, \bar{\lambda}_M))_v$ ;  $\bar{\lambda}_m$  denotes the remaining bandwidth of  
 320  $m$ -th channel.  $\varphi_v$  denotes the number of service requests to be allocated.  $\xi_v$  describes the total bandwidth requests  
 321 of all  $\varphi_v$  cognitive vehicles. In addition, in Eq. (12), the action space is a set composed of whether the number of  
 322  $M$  channels are allocated, in which the action  $a_m$  denotes that the agent allocates the channel  $m$  to a vehicle that  
 323 enters into the priority-based allocation sequence and is ready to be scheduled by the base station currently.

324 A Monte-Carlo search tree consists of nodes and edges. A node  $v$  is a tree node that corresponds to the **MDP**  
 325 **state**  $s_v$ , and the edge connecting a parent node and a child node in the tree represents an action that causes the state  
 326 transition. Each node  $v$  in the tree holds a **node state**, which contains three types of statistics: visit count ( $\bar{N}_v$ ),  
 327 MDP state ( $s_v$ ), and cumulative reward ( $Q_v$ ) received by node  $v$ .

328 The specific search steps are shown in Figure 2.

329 1) Create a root node of the search tree and initialize the node state. Assume that the root node is denoted by  $v$   
 330 and the node state is  $\{\bar{N}_v, s_v, Q_v\}$ .

331 2) Allocate the spectrum resources for vehicles according to the priority order list *Priorityscore\_list* defined  
 332 in Definition 3, and extend the child node while update the node state. Each layer's tree expansion represents the  
 333 spectrum allocation for a vehicle and each allocation process involves many iterations. Take the root node  $v$  in  
 334 Figure 2 as an example. When the channel assignment action of vehicle ID3 is  $a_1$ , the search tree extends down to  
 335 the child node  $v'$  and update the node state through iterative calculation (*i.e.*,  $s_{v'} = f_{ESP}(s_v, a_1)$ ).

336 3) When the tree expansion reaches to the termination condition of iteration (*i.e.*, the second users or the  
 337 available spectrum resources are all allocated), an optimal channel allocation matrix  $A^*$  in the current allocation  
 338 period is returned. For example, assume that when reaching to the node  $v'''$  in Figure 2, the iteration ends. The  
 339 black arrow lines direct an allocation path  $v \rightarrow v' \rightarrow v'' \rightarrow v'''$ . Then the corresponding actions constitute a  
 340 feasible allocation policy set  $\{a_1, a_5, a_1\}$ , which can be converted to a channel allocation matrix  $A_{N \times M}$  as an  
 341 output.

## 342 5.2 DNN-based Environmental State Predictor—ESP

343 Due to the uncertainty of the PUs' spectrum occupancy activities, when the tree is expanded from one node to the  
 344 next in Section , the expansion will be not stable, *i.e.*, given a state and an action, the next state is uncertain. This  
 345 uncertainty is caused by the unknown environment of IoV. Therefore, to limit the expansion scale of the MCTS tree  
 346 horizontally and speed up the search, it is necessary to gradually learn to approach to the real environment of IoV  
 347 when doing spectrum allocation. This section presents an offline environment state predictor (named ESP) based  
 348 on a deep neural network (DNN).

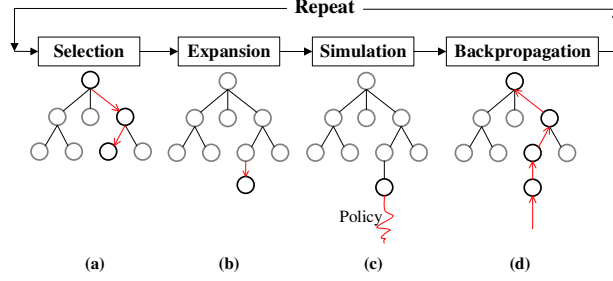


Figure 3: An iterative computation process of Finder-MCTS.

Note that, to obtain the ESP, enough training data are needed. Thus, first during the cold-start phase of Finder-MCTS (*i.e.*, the algorithm just starts running), we do not rely on ESP. This does not affect the channel allocation solution of Finder-MCTS. After a period of time in the cold start phase, our base station can obtain and cumulate large numbers of ‘state-action transition pairs’. Subsequently, we input these ‘state-action transition pairs’ into ESP continuously as the training data to obtain a state transition function  $f_{ESP}$ , which is an offline training process. Once we have the  $f_{ESP}$ , the Finder-MCTS could converge fast due to the reduction in branching. The above training is done by DNN.

The network structure of DNN consists of one input layer, three hidden layer and one output layer. In this paper, we set the learning rate of DNN to 0.05 and the activation function of DNN is the rectified linear unit function (ReLU). To optimize the neural network parameters, we use the mini-batch gradient descent method [22]. In the DNN, the training label is the state  $s_{v'}$ , which is the state of the corresponding expansion child node  $v'$  of node  $v$ . ESP is used to obtain the prediction state  $\hat{s}_{v'}$ . The loss function of ESP is,

$$loss_{ESP} = \frac{1}{B} \sum_B (\|s_{v'} - \hat{s}_{v'}\|_2) \quad (13)$$

where  $B$  represents the batch size of mini-batch gradient descent. In the experiment we set  $B = 64$ , with indicating that 64 samples are selected in each iteration. Notation  $\|\cdot\|_2$  represents the L2 norm. When  $loss_{ESP}$  converges, we let the DNN network parameter  $w_{ESP}$  update.

After we obtain the ESP function, based on the selected action  $a_m$  and MDP state  $s_v$ , ESP can give the MDP state of its expanded node  $\hat{s}_{v'}$ ,

$$\hat{s}_{v'} = f_{ESP}(s_v, a_m | w_{ESP}) \quad (14)$$

### 5.3 Finder-MCTS Algorithm Based on Action Space Pruning and Scenario Simulation

Finder-MCTS requires to execute the following four steps: selection, expansion, simulation, and backpropagation iteratively to complete an computation process, which are shown in Figure 3. In Figure 3, the black circles indicate the nodes involved in each step and the red arrow lines indicate the actions corresponding to each step. In subfigure (c), policy usually refers to the random selection action extended at each step of the simulation process. We usually call step (a) selection and step (b) expansion together as the tree policy. Specifically, the detailed procedures and descriptions are give in the following steps (a)-(d) and in Figure 4.

(a) **Selection.** Each iteration starts from the root node. When the algorithm has to choose to which child node it will descend, it tries to find a good balance between exploitation and exploration. We use the upper confidence bound for tree (UCT) [23] to recursively select child nodes. The selection criterion of the optimal child node is:

$$\arg \max_{v' \in child(v)} \left( \frac{Q_{v'}}{\bar{N}_{v'}} + c \cdot \sqrt{\frac{\ln(\bar{N}_v)}{\bar{N}_{v'}}} \right) \quad (15)$$

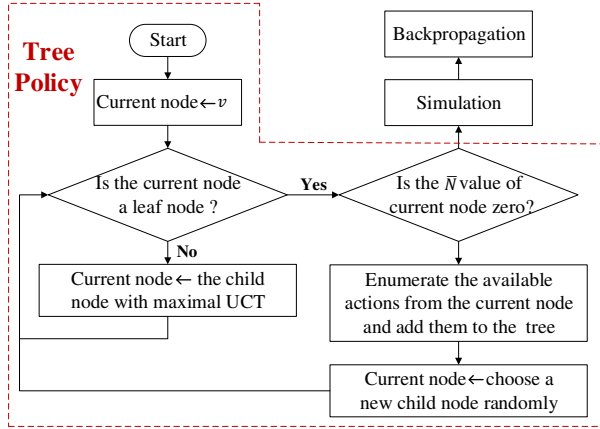


Figure 4: The flow chart of Finder-MCTS.

376 where  $c \geq 0$  is a weight coefficient used to adjust the exploitation and exploration. We set  $c = 0.8$  in the experiment  
 377 through many tests. Notation  $child(v)$  represents the set of child nodes with  $v$  as the parent in the tree.  $\bar{N}_{v'}$  and  
 378  $\bar{N}_v$  represent the total number of times that the child node  $v'$  and its parent node  $v$  have been visited iteratively.  
 379  $Q_{v'}$  represents the cumulative reward obtained by node  $v'$ . Note that, the selected child node should be expandable  
 380 (*i.e.*, have unvisited child node) and represent a non-terminal state. Next, the algorithm treats the child node with  
 381 the largest value of UCT as the current node for the next expansion.

382 (b) **Constraint Oriented Expansion.** Finder-MCTS judges whether the number of visits of the current node is  
 383 0. If visit count  $\bar{N} = 0$ , the algorithm goes to step (c) directly. If the visit count  $\bar{N} \neq 0$ , the algorithm enumerates  
 384 the available actions. However, if it is just a simple enumeration, the number of available actions in the next layer  
 385 is  $M$ . As the tree expands, a huge search tree will be built. The computational complexity grows geometrically  
 386 with the number of SUs to be allocated. Thus, here we give the constraint oriented expansion.

387 In the constraint oriented expansion, we prune the action space according to the constraint conditions defined  
 388 in Section -2) so as to obtain all available actions from the current node. And then add new nodes to expand the  
 389 tree and let the current node be a new child node which is randomly selected after expansion.

390 Specifically, we use  $\mathcal{A}(n, v)$  represent the set of available actions starting from the node  $v$ , which is used for  
 391 the next round of channel allocation for the  $n$ -th SU. That is to say  $\mathcal{A}(n, v)$  is an interference-free action space of  
 392 a SU. The detailed implementation steps of the constraint oriented expansion are described in Algorithm 1.

393 In Algorithm 1, we use three main steps to perform action pruning. First, considering the channel availability,  
 394 we introduce the channel availability matrix  $L$  to prune the set of actions. We map the elements of  $l_{n,m} = 1$   
 395 in the channel availability matrix for vehicle  $n$  to the available action set (Lines 2-6 in Algorithm 1). Second,  
 396 considering that the vehicle  $ID_n$  currently to be allocated should not share the same channel with a vehicle having  
 397 communication interference, we introduce the SU-SU interference matrix  $C$  for the tree pruning. The algorithm  
 398 traverses the elements in the channel allocation matrix  $A$  and makes a judgement on whether  $a_{n',m} = 1$  and  
 399  $c_{n,n',m} = 1$  hold at the same time. If they hold at the same time,  $a_m$  is removed from the action set (Lines 7-15  
 400 in Algorithm 1). Next, in each iteration, the algorithm needs to make a judgement on whether constraint (1) and  
 401 constraints (5-a)  $\sim$  (5-g) hold. If the available channel  $m$  for the vehicle currently to be allocated does not satisfy  
 402 these constraints, action  $a_m$  needs to be removed from the set of actions (Lines 16-20 in Algorithm 1). Finally, If  
 403  $\mathcal{A}(n, v) = \emptyset$ , the algorithm will skip the current allocation and wait for the next round of allocation (Lines 21-23  
 404 in Algorithm 1).

405 (c) **Simulation Based on Different Scenarios.** From the above step (b), we know that if the visit count of the  
 406 current node is zero, we will perform a simulation from the current node (*i.e.*, the newly expanded node, denoted

**Algorithm 1** Constraint Oriented Expansion for Vehicle  $ID_n$ **Input:**

$L$  - channel availability matrix  
 $C$  - SU-SU interference constraint matrix  
 $A$  - channel allocation matrix  
 $\gamma_m$  - the maximum allowable interference level of channel  $m$   
 $\phi_m$  - the available bandwidth of channel  $m$   
 $\delta_{m,k}$  - the maximum allowable interference power of PU  $k$  on channel  $m$

**Output:**

$\mathcal{A}(n, v)$  - the action space/set of vehicle  $ID_n$  under the current node  $v$   
*Function*  $Action(n, v)$

```

1:  $\mathcal{A}(n, v) \leftarrow \emptyset$ 
2: for each  $l_{n,m}$  in the  $n$ -th row of matrix  $L$  do
3:   if  $l_{n,m} = 1$  then
4:      $\mathcal{A}(n, v) \leftarrow a_m$ 
5:   end if
6: end for
7: for each  $c_{n,n',m}$  in  $1 \sim n$  columns of the  $n$ -th row of matrix  $C$  do
8:   for each  $a_{n',m}$  in  $A$  do
9:     if  $c_{n,n',m} = 1$  and  $a_{n',m} = 1$  then
10:      if  $a_m \in \mathcal{A}(n, v)$  then
11:        remove  $a_m$  from  $\mathcal{A}(n, v)$ 
12:      end if
13:    end if
14:  end for
15: end for
16: for each  $a_m$  in  $\mathcal{A}(n, v)$  do
17:   if the available channel  $m$  for the vehicle  $ID_n$  does not satisfy the constraint (1) and constraints (5-a)  $\sim$  (5-g) then
18:     remove  $a_m$  from  $\mathcal{A}(n, v)$ 
19:   end if
20: end for
21: if  $\mathcal{A}(n, v) = \emptyset$  then
22:   the algorithm does not perform the allocation for vehicle  $ID_n$  and waits for the allocation of the next user according to the  $Priorityscore\_list$ 
23: end if

```

407 by  $\tilde{v}$ <sup>1</sup> to the terminal node (denoted by  $\tilde{v}_\Delta$ ). Here, the terminal node refers to the node that the descending arrives  
408 at when the SUs or the available channel resources have been all allocated. Usually, the simulation uses a random  
409 search strategy to generate a reward  $Q_{\tilde{v}_\Delta}$  at the final leaf node  $\tilde{v}_\Delta$ . However, the time-varying property of PUs'  
410 spectrum occupancy activities makes the actual available spectrum resources uncertain. This uncertainty will have  
411 potential impacts on the reward evaluation for the SU to be allocated in IoV.

412 Therefore, in this paper, the duration of network service for a PU (denoted by  $\tau$ ) is included in the simulation  
413 when doing reward evaluation. Reference [24] pointed out that the duration of network service for PU in each  
414 channel obeys a log-normal distribution. The probability density function (PDF) is:

$$f(\tau; \mu, \sigma) = \frac{1}{\tau \sigma \sqrt{2\pi}} e^{-\frac{(\ln \tau - \mu)^2}{2\sigma^2}} \quad (\tau > 0) \quad (16)$$

415 The parameters  $(\mu, \sigma)$  are in milliseconds (ms) and the values used in this paper are (2.47, 1.88) [24].

416 Through random sampling from the above distribution, we can obtain different scenarios of the service durations  
417 for the PUs at each layer in the simulation stage. Each sampling corresponds to a scenario. Since there are infinite  
418 scenarios when sampling, here we sample number of  $\chi$  times at each layer of simulation to control the computation  
419 scale. Thus, a scenario set is formed, denoted by  $\hat{\pi} = \pi^1, \pi^2, \dots, \pi^i, \dots, \pi^\chi$ . In the experiment, we set  $\chi = 9$ . Next,

<sup>1</sup>We use symbol  $\sim$  to label the nodes in the stage of simulation.

we define a stochastic bonus to adjust the reward evaluation according to different service durations, the resource supply and demand situation, and the utilities of SUs.

*Definition 4: Stochastic Bonus*

Assume that the channel  $m$  matches the vehicle  $ID_n$  and the tree expands from node  $\tilde{v}$  to node  $\tilde{v}'$  in the simulation stage. Then, we define a stochastic bonus for node  $\tilde{v}$  as  $\mathbb{E}_{i \in \tilde{\pi}}(H_{n,m}^{\tilde{v}}(i))$ , in which  $\mathbb{E}$  represents the expectation of stochastic bonus obtained by vehicle  $ID_n$  in  $\chi$  scenarios. We have

$$H_{n,m}^{\tilde{v}}(i) = \tanh(Utility_n) \cdot \tau_i^{-1} \cdot \left( \frac{\bar{\lambda}_m}{Count(L_m) - Count(A_m)} \right) \quad (17)$$

where  $\tau_i$  ( $1 \leq i \leq \chi$ ) denotes one of the samplings based on distribution  $f(\tau; \mu, \sigma)$ . The larger the value of  $\tau_i$ , the longer the channel occupied by the PUs in this sampling. It indicates the bonus of vehicle  $ID_n$  when doing allocation will be low. Notation  $Utility_n > 0$  represents the network utility score of vehicle  $ID_n$  (Definition 2), which reflects the communication capability of vehicle  $ID_n$  and is used as a weight coefficient here. We utilize the hyperbolic tangent function  $\tanh(\cdot)$  to normalize the value of  $Utility_n$  to the interval  $[0, 1]$ . When the  $Utility_n$  is large, the weight coefficient is closer to 1, which indicates that the vehicle  $ID_n$  with strong communication ability tends to have high bonus. Besides,  $\frac{\bar{\lambda}_m}{Count(L_m) - Count(A_m)}$  measures the remaining minimum average bandwidth available to vehicle  $ID_n$  currently.  $Count(L_m)$  records the number of elements in the  $m$ -th column with value of 1 in matrix  $L$ . Thus  $Count(L_m) - Count(A_m)$  describes the maximum number of allowable access vehicles on channel  $m$  without considering the interference matrix  $C$  and the available bandwidth  $\phi_m$ .

In summary, if a vehicle with strong communication capability, the PUs with low service durations, and the remaining resources are enough, the stochastic bonus will be high.

Based on the above Eq. (17), we have an adjusted reward  $Q_{\tilde{v}}$  for node  $\tilde{v}$  in the simulation stage,

$$Q_{\tilde{v}} = r_{(n,m)} + \mathbb{E}_{i \in \tilde{\pi}}(H_{n,m}^{\tilde{v}}(i)) \quad (18)$$

where  $r_{n,m}$  refers to the immediate reward that channel  $m$  is allocated to vehicle  $ID_n$  (defined in Eq.(2)). For simplicity, we use notation  $Q_{\tilde{v}}$  with omitting the label of  $n$  and  $m$ .

When the simulation reaches to the terminal node  $\tilde{v}_\Delta$ , we can get the simulation cumulative reward  $Q_{\tilde{v}_\Delta}$  of all nodes on the simulation path from  $\tilde{v}$  to  $\tilde{v}_\Delta$ . We have

$$Q_{\tilde{v}_\Delta} = \sum_{\tilde{v}}^{\tilde{v}_\Delta} \{r_{n,m} + \mathbb{E}_{i \in \tilde{\pi}}(H_{n,m}^{\tilde{v}}(i))\} \quad (19)$$

(d) **Backpropagation.** The aim of backpropagation is to update the empirical information of the prior exploration before the next iteration, which is shown in Figure 5. When an iteration reaches to the terminal node  $\tilde{v}_\Delta$ , according to Eq. (19), we get the simulation cumulative reward  $Q_{\tilde{v}_\Delta}$  for backpropagation.

In this way, the reward of backpropagation can include the reward evaluation of all expanded nodes on the simulation path, with reflecting the overall spectrum allocation performance of simulation in the current iteration. Meanwhile, the algorithm updates the node state on the path from the root to the expanded node according to the following rules:

$$\bar{N}_v \leftarrow \bar{N}_v + 1 \quad (20)$$

$$Q_v \leftarrow Q_v + Q_{\tilde{v}_\Delta} \quad (21)$$

To sum up, we provide the pseudocode of Finder-MCTS in Algorithm 2. The Finder-MCTS algorithm iteratively executes functions such as *Treepolicy*, *Simulation*, and *Backpropagation* to explore different spectrum allocation schemes (i.e.,  $A_m$  in  $\Lambda_{L,C}$ ). It finally finds the optimal spectrum allocation scheme  $A^*$  in the current network.

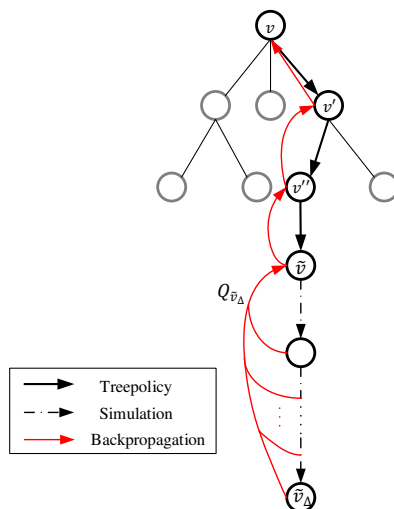


Figure 5: Backpropagation of Finder-MCTS.

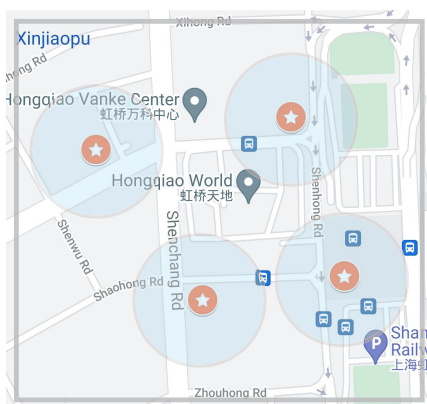


Figure 6: The experimental area imported from OpenStreetMap.

## 455 6. Experimental Results and Analysis

456 In this section, first we give the detailed simulation settings, including the vehicular dataset generation and some  
 457 parameters in our proposed method. Second, we compare Finder-MCTS with other types of methods in terms of  
 458 channel utilization ratio (CUR), average link capacity (ALC), and convergence time. Finally, we test the perfor-  
 459 mance of Finder-MCTS compared with other MCTS algorithms' variations.

### 460 6.1 Simulation Settings

461 Our experiments are done by using the simulation of urban mobility (SUMO) simulator. All the simulations are  
 462 conducted in a PC with Intel Core CPU i9-9820X 3.50GHz processor, 64GB RAM. We export a map of area near  
 463 Pudong Airport in Shanghai from OpenStreetMap, which is shown in Figure 6. The latitude of the experimental area  
 464 is between [31.19177, 31.19742]. The longitude is between [121.31134, 121.31853]. In this area, we randomly  
 465 select four base stations (depicted by red star marks). The locations of these base stations and different communi-  
 466 cation radius are listed in TABLE 1. Each base station can observe the traffic flows and obtain the passing vehicles'  
 467 information, including the vehicle ID, location, speed, timestamp and acceleration.

Table 1: Information of the four base stations.

Name	Latitude	Longitude	Communication Radius
BS1	31.19554	121.31274	500m
BS2	31.19604	121.31619	500m
BS3	31.19327	121.31462	500m
BS4	31.19363	121.31713	500m

Table 2: Parameters used in SUMO.

Parameters	Car	Bus	Truck
the maximum speed	15(m/s)	13(m/s)	10(m/s)
the minimum speed	1(m/s)	1(m/s)	1(m/s)
the minimum gap between vehicles	2.5(m)	2.5(m)	2.5(m)
the maximum acceleration	3(m/s <sup>2</sup> )	1.5(m/s <sup>2</sup> )	1.5(m/s <sup>2</sup> )
the maximum deceleration	7.5(m/s <sup>2</sup> )	4(m/s <sup>2</sup> )	4(m/s <sup>2</sup> )
the maximum deceleration for emergency breaking	9(m/s)	7(m/s)	7(m/s)

468 Assume that each base station has  $M = 10$  available spectrum channels. The bandwidth of each channel is set  
 469 to 20MHz. We import 100 cognitive vehicles into the simulation scene. Each vehicle randomly proposes a service  
 470 request to the base station with probability of 50% at each allocation time window. Suppose that the duration of  
 471 network service for each vehicle is equal to the allocation time window. In SUMO, we set the parameters for  
 472 the different types of vehicles in TABLE 2. Compared with the moving vehicle, a PU can be regarded as a static  
 473 point in the experiment. We set a total of  $K = 30$  fixed points as PUs under the four base stations. Each PU  
 474 randomly occupies a part of the communication bandwidth (MHz), which subjects to  $U[1, 10]$  uniform distribution.  
 475 At each allocation time window, we randomly let 70% PUs occupy the nearest base station's available channels.  
 476 The duration of network service for a PU is chosen according to Eq.(16). The spectrum demand of each SU  $n$  are  
 477 randomly selected in [1,3] MHz. The maximum allowable interference level on channel  $m$   $\gamma_m$  is  $-114dBm$ . The  
 478 level of background noise on channel  $m$   $N_m$  is 1dB. The minimum transmission power and maximum transmission  
 479 power are  $\bar{P}_{m,n}^{min} = 20dBm$  and  $\bar{P}_{m,n}^{max} = 25dBm$  respectively. The maximum allowable interference power of PU  
 480  $k$  on channel  $m$   $\delta_{m,k}$  is 5dB.

481 In the experiment, the protection radius of a PU ( $\tilde{R}_{k,m}$ ) is set to 100m. We let the the transmit power level of SUs  
 482 be generated from the set  $[20dBm, 21dBm, 22dBm, 23dBm, 24dBm, 25dBm]$ . Thus, the interference radius of a  
 483 SU ( $\bar{R}_{n,m}$ ) corresponding to the above power levels are 100m, 150m, 200m, 250m, 300m and 350m. The transmit  
 484 power of a base station is set to 46dBm. For simplicity, assume that the transmit power is equal to the transmission  
 485 power and let the channel gain in the wireless space be constituted by the path loss. We define the path loss between  
 486 SU  $n$  and base station  $j$  as  $PL(n, j) = 34 + 40lg(d_{nj})$ , where  $d_{nj}$  denotes the Euclidean distance between SU  $n$   
 487 and base station  $j$ . Besides, we define the path loss between SU  $n$  and PU  $k$  as  $PL(n, k) = 40 + 24.4lg(d_{nk})$ [25],  
 488 where  $d_{nk}$  denotes the Euclidean distance between SU  $n$  and PU  $k$ . The received signal power level is given by the  
 489 product of the transmit power and the channel gain. Thus the parameters  $P_{m,n,k}$ ,  $P_{m,k,n}$ , and  $P_{m,n}$  can be obtained  
 490 through the above calculations.

## 491 6.2 Comparison with Other Types of Methods

492 Under the same simulation settings, we compare our Finder-MCTS with three other algorithms, *i.e.*, the game  
 493 theory-based method [6], particle swarm optimization-based (PSO-based) method [11], and DQN-based method  
 494 [26], in terms of channel utilization ratio, convergence time and average link capacity of SUs.



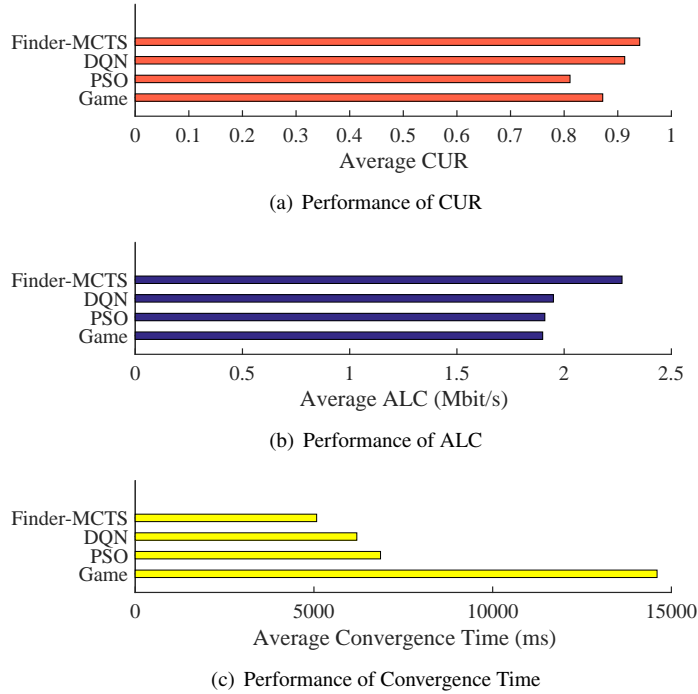


Figure 7: Comparison with three other methods in terms of average CUR, ALC and convergence time.

495 The channel utilization ratio (CUR) refers to the occupancy ratio of the available spectrum resources in the  
 496 current base station. Besides, the average link capacity (ALC) is defined as follows:

$$ALC = \frac{1}{N} \cdot \sum_{n=1}^N \sum_{m=1}^M (a_{n,m} \cdot r_{n,m}) \quad (22)$$

497 If a method has high CUR, high ALC and low convergence time, it means that the method can not only make  
 498 full use of the spectrum resources, but also can enable SUs to obtain better communication service quality quickly.

499 First, after the simulations are all done in the four base stations, we compare the average CUR, ALC, and  
 500 convergence time of the proposed Finder-MCTS with three other methods, shown in Figure 7. From the average  
 501 CUR performance in Figure 7 (a), we can see that Finder-MCTS performs the best, the second-best is DQN-based  
 502 method, and the worst is PSO-based method. From the average ALC performance in Figure 7 (b), we can see that  
 503 Finder-MCTS performs the best, the second-best best is DQN-based method, and the worst is game theory-based  
 504 method. From the average convergence time performance in Figure 7 (c), we can see that Finder-MCTS performs  
 505 the best, the second-best best is DQN-based method, and the worst is game theory-based method.

506 Based on the above results, we give the following analysis. Because the convergence of the Nash equilibrium  
 507 solution is negatively related to the size of the problem, the game theory-based method's convergence performance  
 508 is poor. When the game theory-based method reaches convergence, the CUR performance of the system can  
 509 be approximately optimal, however the equilibrium of the multi-user game makes the ALC value relatively low.  
 510 Besides, the PSO-based method is easy to fall into the local optimal solution, its average CUR and average ALC  
 511 perform relatively poor. Since the complicated parameters' setting of PSO, its average convergence time becomes  
 512 longer as the scale of the problem becomes larger. Moreover, after the exploration of actions through reinforcement  
 513 learning, DQN-based method can obtain a higher quality spectrum allocation solution, and the performance of  
 514 average CUR and ALC is second only to Finder-MCTS. However, the convergence time of DQN-based method  
 515 is higher than Finder-MCTS due to the long-term exploration and value updating, although enough experience

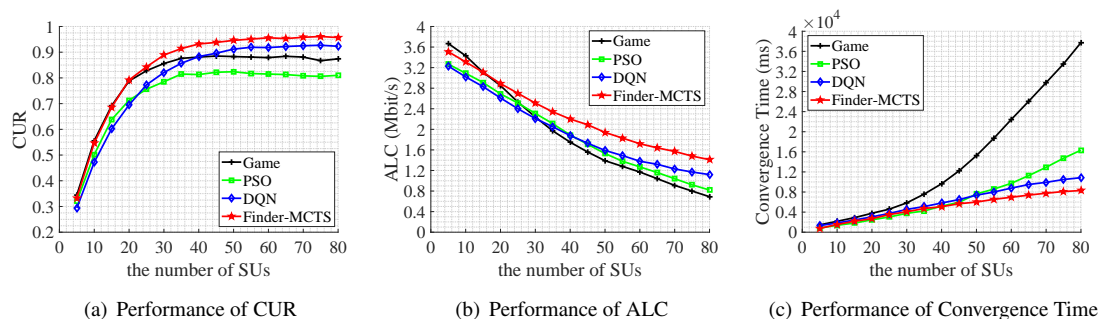


Figure 8: Performance comparison with varying the number of SUs.

516 information learned through online learning can speed up the convergence time of DQN to some extent. By contrast,  
 517 Finder-MCTS based on offline training and online learning has an average 36.47% improvement in convergence  
 518 time than other methods. In terms of ALC, Finder-MCTS has an average advantage of 18.24% over other methods.  
 519 At the same time, the channel utilization of Finder-MCTS is 9.00% higher than other methods on average.

520 Second, since the number of SUs in the coverage area of each base station is time-varying, it is necessary to  
 521 observe the performance changes under different SUs' scales. The results are shown in Figure 8. Here, notice that  
 522 in Figure 8, each depicted point in the curve is an averaged value statistically. For example, as to the results that dis-  
 523 tribute in the scale interval  $(p_1, p_2]$  of x-axis, we average these results and depict the averaged value corresponding  
 524 to point  $p_2$ .

525 Figure 8 (a) shows the relationship between the number of SUs and CUR. In general, as the number of SUs  
 526 increases, the CUR curve increases until it gradually converges. In addition, we find that when the number of SUs  
 527 is small, the game theory can give a solution with high CUR. However, with the increase of SUs, Finder-MCTS  
 528 and DQN-based method show obvious advantages in resource utilization. The reason behind that is when the scale  
 529 of SU becomes large, the combination of historical experiences and online exploration can greatly improve the  
 530 quality of the solution. In contrast, the game theory-based equilibrium quality for large-scale SU problem has  
 531 declined. Also, the PSO-based method often converges to a local optimal solution and its CUR performance cannot  
 532 be guaranteed.

533 Figure 8 (b) depicts the relationship between the number of SUs and ALC. It is obvious that as the number of  
 534 SUs increases, the ALC value decreases since the available spectrum resources of the base station side are limited.  
 535 Besides, we find that when the number of SUs is small, the game-based method shows a good performance in  
 536 ALC. However, as the number of SUs increases, Finder-MCTS shows an obvious advantage. This because when  
 537 the scale of SUs becomes large, finding an optimal solution is hard for the game-based method. Moreover, since  
 538 the PSO-based method is hard to reach the global convergence, the ALC performance is relatively low with the  
 539 number of SUs increasing.

540 Figure 8 (c) shows the simulation results of the relationship between the number of SUs and the convergence  
 541 time. First, we can see that the convergence time of game theory-based and PSO-based method shows an obvi-  
 542 ous growth trend as the number of SUs increases, while the convergence time based on DQN and Finder-MCTS  
 543 rises moderately. The main reason is that the Finder-MCTS and DQN-based methods gradually fit the channel  
 544 state model after continuous learning, thereby greatly improving the search efficiency. The convergence time of  
 545 Finder-MCTS is reduced by 65.23% and 18.85% compared with the game theory-based method and the PSO-based  
 546 method. In the long run, Finder-MCTS shows a short and gentle convergence time performance in the dynamic  
 547 environment.

548 All above phenomena verify the advantage of Finder-MCTS in solving spectrum allocation in IoV. Finder-  
 549 MCTS can effectively complete the rapid learning of the approximate optimal allocation solution in a time-varying  
 550 environment, which greatly improves the available spectrum utilization ratio of the current base station system.

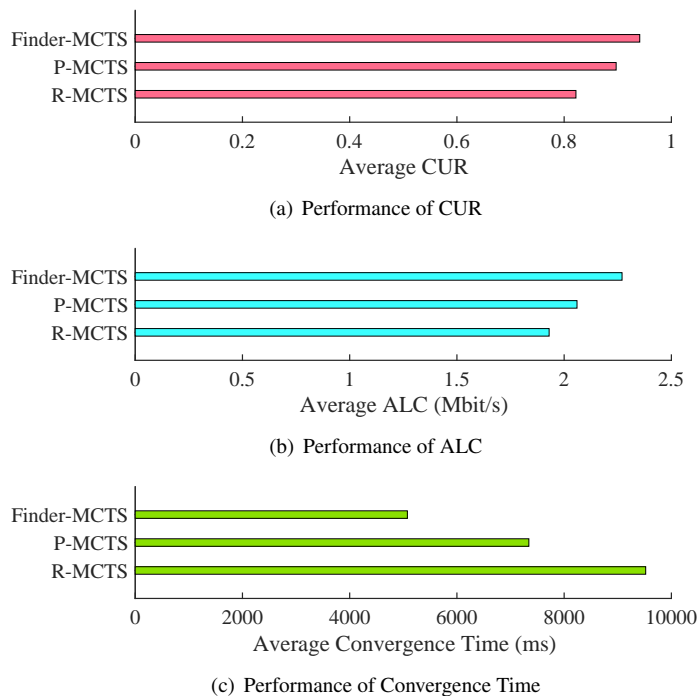


Figure 9: Comparison with two types of MCTS algorithms' variations in terms of average CUR, ALC and convergence time.

### 6.3 Comparison with Other MCTS Algorithms' Variations

In this part, we compare Finder-MCTS with other MCTS algorithms' variations. We show why we consider the priority mechanism and simulation under different scenarios.

We set two basic types of MCTS-based spectrum allocation modes: random-order-based allocation mode and priority-based allocation mode, which are called as R-MCTS and P-MCTS respectively. In R-MCTS, compared with Finder-MCTS, both priority and the uncertainty of PUs' service durations are not taken into consideration. In P-MCTS, compared with Finder-MCTS, only the uncertainty of PUs' service durations is not taken into consideration. The simulation results are shown in Figure 9. We can see that Finder-MCTS performs the best, the second-best is P-MCTS, and the worst is R-MCTS. According to the above results, we give the following analysis.

From Figure 9 (a), we can see that the CUR performance of P-MCTS is superior than R-MCTS. This gap illustrates that the introduction of priority evaluation will improve the ratio of the spectrum utilization (about 9.12% increase). Meanwhile, Finder-MCTS has the best CUR performance. In the long run, the service duration  $\tau$  of PU on each channel will give each allocated SU differentiated stochastic bonus. Hence, based on the uncertainty of the channel state occupied by the PUs, we introduce the factor  $\tau$  that affects the supply-demand ratio of spectrum resources into the reward evaluation during each expansion step of the simulation process. We learn about Finder-MCTS is better (about 4.08% increase) than P-MCTS on ALC. Hence, we can conclude that the optimization of the stochastic simulation process contribute to improve spectrum usage efficiency of CR-IoV from a global perspective.

Figure 9 (b) depicts the different performances of the three methods in ALC performance. With the help of priority evaluation, P-MCTS has increased by 6.73% compared with R-MCTS. The ALC performance of Finder-MCTS has increased by 10.19% compared with P-MCTS by evaluating the uncertainty of PUs' service durations.

Figure 9 (c) shows the average convergence time of three methods. Owing to the priority evaluation, P-MCTS has a 22.89% advantage than R-MCTS. This characterizes the positive impact of the differentiation priority evaluation on the algorithm convergence time. Secondly, under the same setting, with the help of reduction of action

574 space in each descending layer, Finder-MCTS achieves a faster convergence speed (about 46.69% increase and  
575 30.86% increase) than R-MCTS and P-MCTS.

## 576 7. Conclusion

577 In this paper, we investigate the spectrum allocation in CR-IoV by modeling a optimization problem to maximize  
578 the link capacity of vehicle users. What's more, we propose a method named Finder-MCTS to solve the optimiza-  
579 tion problem. We show that Finder-MCTS can learn to adapt and update allocation strategy for transmission under  
580 dynamic network environment. The experimental results show that Finder-MCTS is more efficient in convergence  
581 speed, and it achieves good performance gain in spectrum utilization and link capacity compared with other pop-  
582 ular strategies, especially when the number of vehicle users becomes more. Besides, we have also confirmed the  
583 effectiveness of priority evaluation and uncertainty evaluation of the PUs' service durations by comparing with two  
584 variations of MCTS. In the future work, we will further study the cooperative spectrum allocation problem of IoV  
585 under a complex scenario with space/air/ground communications and networking.

## 586 Data Availability

587 The data generation method has been introduced in section 6.1. The data can be obtained according to the configu-  
588 ration. We also make data available on request through sending email to the authors.

## 589 Conflicts of Interest

590 The authors declare that they have no conflicts of interest.

## 591 Funding

592 This work was supported in part by the National Natural Science Foundation of China under Grant 61972080, in  
593 part by the Shanghai Rising-Star Program under Grant 19QA1400300.

## 594 References

- 595 [1] Jian Yang and Hangsheng Zhao, "Enhanced throughput of cognitive radio networks by imperfect spectrum prediction",  
596 *IEEE Communications Letters*, vol. 19, no. 10, pp. 1738–1741, 2015.
- 597 [2] Mohammad Yousefvand, Nirwan Ansari and Siavash Khorsandi, "Maximizing network capacity of cognitive radio net-  
598 works by capacity-aware spectrum allocation", *IEEE Transactions on Wireless Communications*, vol. 14, no. 9, pp. 5058–  
599 5067, 2015.
- 600 [3] John Martinovic, Eduard Jorswieck, Guntram Scheithauer and Andreas Fischer, "Integer linear programming formulations  
601 for cognitive radio resource allocation", *IEEE Wireless Communications Letters*, vol. 6, no. 4, pp. 494–497, 2017.
- 602 [4] Zhijun Teng, Luying Xie, Haolei Chen, Lixin Teng and Hongbiao Li, "Application research of game theory in cognitive  
603 radio spectrum allocation", *Wireless Networks*, vol. 25, no. 7, pp. 4275–4286, 2019.
- 604 [5] Changyan Yi and Jun Cai, "Two-stage spectrum sharing with combinatorial auction and stackelberg game in recall-based  
605 cognitive radio networks", *IEEE Transactions on Communications*, vol. 62, no. 11, pp. 3740–3752, 2014.
- 606 [6] Xiaozhu Liu, Rongbo Zhu, Brian Jalaian and Yongli Sun, "Dynamic spectrum access algorithm based on game theory in  
607 cognitive radio networks", *Mobile Networks and Applications*, vol. 20, no. 6, pp. 817–827, 2015.

- 608 [7] Zhijin Zhao, Zhen Peng, Shilian Zheng and Junna Shang, “Cognitive radio spectrum allocation using evolutionary algo-  
609 rithms”, *IEEE Transactions on Wireless Communications*, vol. 8, no. 9, pp. 4421–4425, 2009.
- 610 [8] Albert Y.S. Lam, Victor O.K. Li and James J.Q. Yu, “Power-controlled cognitive radio spectrum allocation with chemical  
611 reaction optimization”, *IEEE Transactions on Wireless Communications*, vol. 12, no. 7, pp. 3180–3190, 2013.
- 612 [9] Piyush Bhardwaj, Ankita Panwar, Onur Ozdemir et al., “Enhanced dynamic spectrum access in multiband cognitive radio  
613 networks via optimized resource allocation”, *IEEE Transactions on Wireless Communications*, vol. 15, no. 12, pp. 8093–  
614 8106, 2016.
- 615 [10] Ruining Zhang, Xuemei Jiang and Ruifang Li, “Improved decomposition-based multi-objective cuckoo search algorithm  
616 for spectrum allocation in cognitive vehicular network”, *Physical Communication*, vol. 34, pp. 301–309, 2019.
- 617 [11] Quan Liu, Hongwei Niu, Wenjun Xu and Duzhong Zhang, “A service-oriented spectrum allocation algorithm using en-  
618 hanced pso for cognitive wireless networks”, *Physical Communication*, vol. 74, no. 12, pp. 81–91, 2019.
- 619 [12] Qian Huang, Xianzhong Xie, Hong Tang et al., “Machine-learning-based cognitive spectrum assignment for 5g urllc  
620 applications”, *IEEE Network*, vol. 33, no. 4, pp. 30–35, 2019.
- 621 [13] Feng Li, Dongxiao Yu, Huan Yang et al., “Multi-armed-bandit-based spectrum scheduling algorithms in wireless networks:  
622 A survey”, *IEEE Wireless Communications*, vol. 27, no. 1, pp. 24–30, 2020.
- 623 [14] Yi Zhang, Wee Peng Tay, Kwok Hung Li, Moez Essegir and Dominique Gaïti, “Learning temporal&spatial spectrum  
624 reuse”, *IEEE Transactions on Communications*, vol. 64, no. 7, pp. 3092–3103, 2016.
- 625 [15] Oshri Naparstek and Kobi Cohen, “Deep multi-user reinforcement learning for distributed dynamic spectrum access”,  
626 *IEEE Transactions on Wireless Communications*, vol. 18, no. 1, pp. 310–323, 2019.
- 627 [16] Zhiming Hu, James Tu and Baochun Li, “Spear: Optimized dependency-aware task scheduling with deep reinforcement  
628 learning”, in *Proceedings of the IEEE ICDCS*, pp. 2037–2046, Dallas, TX, USA, Jul. 2019.
- 629 [17] Hang Shuai and Haibo He, “Online scheduling of a residential microgrid via monte-carlo tree search and a learned model”,  
630 *IEEE Transactions on Smart Grid*, vol. 12, no. 2, pp. 1073–1087, 2021.
- 631 [18] Elias Z. Tragos, Sherali Zeadally, Alexandros G. Fragkiadakis and Vasilios A. Siris, “Spectrum assignment in cognitive  
632 radio networks: A comprehensive survey”, *IEEE Communications Surveys Tutorials*, vol. 15, no. 3, pp. 1108–1135, 2013.
- 633 [19] R Sutton and A Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.
- 634 [20] Hailan Yang, Shengze Li, Xinhai Xu et al., “Efficient searching with mcts and imitation learning: A case study in pom-  
635 merman”, *IEEE Access*, vol. 9, pp. 48851–48859, 2021.
- 636 [21] Y. Shang, W. Wu, J. Guo and J. Liao, “Stochastic maintenance schedules of active distribution networks based on monte-  
637 carlo tree search”, *IEEE Transactions on Power Systems*, vol. 35, no. 5, pp. 3940–3952, 2020.
- 638 [22] Mu Li, Tong Zhang, Yuqiang Chen and Alexander J. Smola, “Efficient mini-batch training for stochastic optimization”, in  
639 *Proceedings of the ACM SIGKDD*, pp. 661–670, New York, NY, USA, Aug. 2014.
- 640 [23] Levente Kocsis and Csaba Szepesvári, “Bandit based monte-carlo planning”, in *Proceedings of the ECML*, pp. 282–293,  
641 Berlin, Germany, Sep. 2006.
- 642 [24] Ki Won Sung, Seong-Lyun Kim and Jens Zander, “Temporal spectrum sharing based on primary user activity prediction”,  
643 *IEEE Transactions on Wireless Communications*, vol. 9, no. 12, pp. 3848–3855, 2010.
- 644 [25] Ibrahim Rashdan, Fabian de Ponte Mízzler, Wei Wang, Martin Schmidhammer and Stephan Sand, “Vehicle-to-pedestrian  
645 channel characterization: Wideband measurement campaign and first results”, in *Proceedings of the EuCAP*, pp. 1–5,  
646 London, UK, Apr. 2018.
- 647 [26] Wanlu Lei, Yu Ye and Ming Xiao, “Deep reinforcement learning-based spectrum allocation in integrated access and  
648 backhaul networks”, *IEEE Transactions on Cognitive Communications and Networking*, vol. 6, no. 3, pp. 970–979, 2020.

**Algorithm 2** Finder-MCTS**Input:***Priorityscore\_list***Output:**optimal channel allocation matrix  $A^*$ *Function Finder – MCTS*( $v, Priorityscore\_list$ )

```

1: load network  $f_{ESP}$ 
2: create root node  $v$  with state  $s_v$ 
3: create channel allocation buffer  $\Lambda_{L,C}$ 
4: while node  $v$  is a terminal node do
5:   initialize a matrix  $A_{N \times M}$  with all elements equaling to 0
6:    $\tilde{v} \leftarrow Treepolicy(v)$ 
7:    $Q_{\tilde{v}\Delta} \leftarrow Simulation(s_{\tilde{v}}, \tilde{v})$ 
8:   if  $a_m = 1$  for vehicle  $ID_n$  then
9:      $a_{n,m} = 1$ 
10:  else
11:     $a_{n,m} = 0$ 
12:  end if
13:  update and put  $A_{N \times M}$  in  $\Lambda_{L,C}$ 
14:  Backpropagation( $v, Q_{\tilde{v}\Delta}$ )
15: end while
16: return  $A^* = \underset{A_{N \times M} \in \Lambda_{L,C}}{argmax} \{U(A_{N \times M}, R)\}$ 
Function Treepolicy( $v$ )
17: while  $v$  is nonterminal do
18:  if  $v$  is not a leaf node then
19:     $v' \leftarrow Bestchild(v)$ 
20:    Treepolicy( $v'$ )
21:  else
22:    if  $\bar{N}_v = 0$  then
23:       $\tilde{v} \leftarrow v$ 
24:    else
25:      Expand( $v$ )
26:    end if
27:  end if
28: end while
Function Bestchild( $v$ )
29: return  $\underset{v' \in child(v)}{argmax} (\frac{Q_{v'}}{\bar{N}_{v'}} + c \cdot \sqrt{\frac{\ln(\bar{N}_v)}{\bar{N}_{v'}}})$ 
Function Expand( $v$ )
30: execute Action( $n, v$ )
31: choose  $a_m \in \mathcal{A}(n, v)$  randomly
32: generate a new child  $v'$  of node  $v$ 
33: initialize  $Q_{v'} = 0$ 
34:  $s_{v'} = f_{ESP}(s_v, a_m)$ 
35: Treepolicy( $v'$ )
Function Simulation( $\tilde{v}$ )
36: initialize  $i = 0, Q_{\tilde{v}} = 0$ 
37: while  $\tilde{v}$  is not a terminal node  $\tilde{v}_\Delta$  do
38:  choose  $a_m \in \mathcal{A}(n, \tilde{v})$  randomly
39:   $s_{\tilde{v}'} \leftarrow f(s_{\tilde{v}}, a_m), \tilde{v}' \leftarrow \tilde{v}$ 
40:  calculate  $r_{n,m}$  according to Eq. (2)
41:   $Q_{\tilde{v}'} \leftarrow Q_{\tilde{v}} + r_{n,m} + Bonus$  (Bonus is calculated based on Eq. (17)-(19))
42:   $i \leftarrow i + 1$ 
43: end while
44: return  $Q_{\tilde{v}\Delta}$  when node  $\tilde{v}$  reaching to the terminal node  $\tilde{v}_\Delta$ 
Function Backpropagation( $v, Q_{\tilde{v}\Delta}$ )
45: while node  $v$  is not null do
46:   $\bar{N}_v \leftarrow \bar{N}_v + 1, Q_v \leftarrow Q_v + Q_{\tilde{v}\Delta}$ 
47:   $v \leftarrow parent\ of\ v$ 
48: end while

```