Multi-Agent Reinforcement Learning based Adaptive Parameter Optimization for Semi-Persistent Scheduling in C-V2X Mode 4

Pengcheng Fan¹, Xinyu Chen¹, Jing Zhao^{1*}, Ning Lu¹, Ping Wang¹

¹ School of Software Technology, Dalian University of Technology, Dalian, 116620, China {fanpc, zhangjiaming, luning0408, 2019wangping}@mail.dlut.edu.cn, zhaoj9988@dlut.edu.cn

Abstract—The Third Generation Partnership Project (3GPP) has standardized cellular vehicle-to-everything (C-V2X) Mode 4 to facilitate direct communication between intelligent connected vehicles. In Mode 4, vehicles autonomously reserve and select wireless spectrum resources through sensing-based semipersistent scheduling (SPS). However, the half-duplex transmissions and hidden terminal problems in SPS could degrade the quality of services (QoS), possibly making it hard to meet the requirements for basic safety services. To support the SPS algorithm's performance in highly congested conditions, this paper presents QMIX-SPS, an adaptive parameter optimization methodology. It also proposes a new performance metric, the signal-to-interference-plus-noise ratio (SINR) achievement rate, as a means of evaluating communication network effectiveness. We developed a multi-agent vehicular communication framework in which parameters of the SPS, including transmission power, resource reservation probabilities, resource reservation counteroffset steps, and candidate resource ratios, are periodically adjusted under the guidance of the QMIX reinforcement learning algorithm. Our proposed resource allocation algorithm utilizes the mechanism of value decomposition to solve the reward allocation problem when competition and collaboration coexist in a V2X environment. Simulation results show that QMIX-SPS outperforms the baseline algorithm. Furthermore, the algorithm has excellent stability flexibility, and compatibility with the SPS.

Index Terms—C-V2X Mode 4, spectrum resource allocation, semi-persistent scheduling, multi-agent reinforcement learning

I. INTRODUCTION

The Internet of Vehicles (IoV) has become an essential component in enhancing road traffic systems' safety, efficiency, and convenience due to the development of Intelligent Transport Systems (ITS). In recent years, vehicle-to-everything (V2X) communications have gained the interest of both business and academics, emerging as a key component of vehicular technology.

Two main technologies facilitate vehicle network communication: Cellular Vehicle-to-Everything (C-V2X) and Dedicated Short-Range Communication (DSRC) [1]. While C-V2X was introduced by the Third Generation Partnership Project (3GPP) in its Release 14 [2], DSRC is based on the IEEE 802.11p standard. Comparing C-V2X to DSRC, greater scalability, higher Quality of Service (QoS), and a larger communication coverage range are provided by utilizing LTE and 5G technology. Mode 3 and Mode 4 are the two radio resource allocations that C-V2X provides for vehicle-to-vehicle (V2V) direct communication. Mode 4 enables autonomous resource selection by vehicles via the PC5 interface, distributing radio resources according to sensing-based semi-persistent scheduling (SPS). Mode 3 entails centralized resource scheduling through the base station's UU interface. The paper will concentrate on the study of the SPS for vehicle-to-vehicle (V2V) communication within C-V2X Mode 4.

The performance of the V2X resource allocation algorithms is extremely important for achieving reliable direct communication and broadcasting safety messages between adjacent vehicles. In previous work [3-6], people have studied the performance of the SPS algorithm through analytical models and simulation models and found that appropriately adjusting SPS parameters can provide better performance. Dayal et al. [7] proposed an adaptive SPS scheme to adjust the resource reservation interval under different vehicle traffic scenarios. This reduces interference between adjacent vehicles and increases the effective communication distance. Based on the analytical model. Gu et al. [8] optimized parameters such as the SPS resource reservation interval to reduce the congestion level of the channel. Although these methods have optimized SPS parameters to a certain extent, there are still problems with half-duplex transmission and hidden terminals, so some research efforts have turned to improving the SPS process. Kim et al.[9] proposed a method based on intelligent part sensing that enhances the sensing capability of SPS by minimizing the number of blind decodes. Wang et al. [10] proposed an in-platoon collaborative sensing method to solve the hidden terminal problem in the platoon scenario.

With the introduction of more advanced V2X applications, traditional optimization schemes cannot meet the diverse performance requirements of resource allocation algorithms, and some studies have shown great potential by introducing reinforcement learning-based resource allocation schemes through Markov modeling. Liang et al. [11] innovatively introduced a multi-intelligent resource allocation method to solve the different QoS requirements of V2V and V2I links. Parvini et al. [12] proposed a task reward decomposition mechanism to further improve the performance of the resource allocation algorithm based on deep reinforcement learning (DRL) in the formation scenario.

Based on the research of related work, this paper proposes a DRL-based SPS algorithm that integrates the value decomposition network called QMIX-SPS [13]. We provide a novel met-

ric called the Signal-to-Interference-plus-Noise Ratio (SINR) achievement rate to carry out a comprehensive optimization of the network performance.QMIX-SPS can guide vehicles to select the parameters of SPS periodically to adapt to the rapidly changing traffic conditions and network topology. We use an architecture combining distributed execution and centralized training to handle the value assignment problem in a multi-agent environment. The main contributions of our work can be summarized in three aspects:

- For V2V links to meet the reliability and timeliness requirements necessary to provide essential security services for the Internet of Vehicles (IoV), this work develops a utility index, or SINR achievement rate, based on transmission delay and packet transmission interruption rate. The utility index considers the impact of delay and stability on service quality in addition to the probability of successful transmission.
- This paper presents a novel attempt at simultaneously optimizing SPS parameters through reinforcement learning techniques. In the case of intense channel competition, the establishment of effective communication and cooperation between vehicles can successfully mitigate resource conflicts and improve the overall communication performance of the network.
- This paper is the initial study that investigates the compatibility of C-V2X resource allocation. We fully considered the impact of introducing QMIX-SPS in the SPS environment. The simulation results show that the QMIX-SPS we proposed has almost no impact on the original SPS environment and also improves the stability and reliability of the overall network.

The subsequent sections of this paper are organized in the following manner. Section II presents our system model. Section III introduces the problem description. Section IV outlines the proposed QMIX-SPS algorithm. Section V evaluates the proposed resource allocation algorithm through simulation. Finally, Section VI provides a summary of the paper.

II. SYSTEM MODEL

Orthogonal Frequency Division Multiplexing (OFDM) technology is used in the Internet of Vehicles. Its principle is to convert selected channels in the frequency domain into parallel flat channels on multiple subcarriers. Several consecutive subcarriers are divided into a spectrum subchannel, assuming that the channel fading is approximately the same within a subchannel and is independent of different subchannels. The available spectrum bandwidth of V2V direct communication is w MHz, which is divided into k subchannels, and it is assumed that different subchannels do not interfere with each other. From the time domain, every 100 ms is a broadcast period. Within the same broadcast period, each vehicle can only reserve one resource block at the same time. The vehicle needs to use the adaptive strategy π to select spectrum resources from the spectrum resource pool to transmit periodic broadcast messages to surrounding vehicles. The transmission power of vehicle *i* in the *t*th period is P_t^i , and the selected spectrum resource is rs_t^i . The transmission gain between vehicles changes dynamically as the vehicle moves. Since the single period time is short, assuming the transmission gain is the same in a single period, the transmission gain of vehicle *i* and vehicle *j* in the *t*th period is:

$$g_t^{i,j} = PL_t^{i,j} + SF \quad (i \neq j) \tag{1}$$

 $PL_t^{i,j}$ and SF are path loss and shadow fading respectively. vehicle j when receiving the signal from vehicle i, the interference from other vehicles is described as:

$$I_t^{(j,i)} = \sum_{i^* \neq j, i^* \neq i} B_t^{(j,i^*)} P_t^{i^*} g_t^{(j,i^*)}$$
(2)

where $B_t^{(j,i^*)}$ is 1 if and only if rs_t^i and $rs_t^{i^*}$ are the same, otherwise it is 0. Then the description of SINR is:

$$\gamma_t^{(j,i)} = \frac{P_t^i g_t^{(i,j)}}{\sigma^2 + I_t^{(j,i)}}$$
(3)

where the numerator term represents the effective received power of vehicle j to vehicle i, and σ^2 is the noise power. According to Shannon's second theorem, the transmission rate of vehicle j receiving vehicle i is $T_t^{(j,i)}$, which is given by:

$$T_t^{(j,i)} = \frac{w}{k} log_2(1 + \gamma_t^{(j,i)})$$
(4)

III. PROBLEM FORMULATION

Various application services in C-V2X have distinct QoS requirements. The V2V direct link is primarily responsible for transmitting safety application messages, such as the CAM information that the vehicle periodically broadcasts, and contains the most important basic vehicle data, including position, speed, direction, and other data. Such messages have a great impact on the security of intelligent network-connected vehicles. This link needs to ensure the stability and timeliness of data packet transmission. We choose to use the packet transmission interruption rate to measure the stability of data packet delivery and the sending delay to measure the timeliness.

Based on [11], it can be obtained that under the constraint of the upper limit of transmission interruption probability p_o , the SINR of data packet transmission needs to satisfy the following formula:

$$\gamma_t^{(j,i)} \ge \frac{\gamma_{min}}{\ln\left(\frac{1}{1-p_o}\right)} \tag{5}$$

where γ_m represents the SINR threshold transmission delay, which is defined as follows:

$$\frac{Z}{T_t^{(j,i)}} \le t d_{max} \tag{6}$$

In (6), Z represents the data packet size, and td_{max} represents the upper limit of the sending delay required by the security application. Combining (5) and (6), we can simplify the service quality requirements into constraints on SINR:

$$\gamma_t^{(j,i)} \ge \max\left(\frac{\gamma_m}{\ln\left(\frac{1}{1-p_o}\right)}, 2^{\frac{Z \cdot k}{td_{max} \cdot w}} - 1\right)$$
(7)

To maximize $\gamma_t^{((j,i)}$ for a single vehicle, one can enhance the transmission power P_t^i . However, this will cause interference with other vehicles that select the same spectrum resources, causing the SINR of these vehicles to decrease. A strategy that is better for some vehicles may not necessarily be better overall. Therefore, the understanding of this problem should be considered from a global perspective. This paper aims to make as much data transmission as possible satisfy the constraints of (7). We design the global SINR achievement rate as an evaluation index to measure strategy performance:

$$PAS = \frac{N_a}{N_t} \tag{8}$$

where N_a represents the number of packet transmissions whose SINR satisfies the constraint condition of (7), and N_t represents the total number of transmissions. This paper hopes to maximize *PAS* by dynamically and adaptively adjusting the parameters and transmit power of the SPS algorithm. Therefore, we have the following formal description of this problem:

$$s.t. \begin{cases} 0 \le H \le 1\\ 0 \le F \le 1\\ C \in \{0, 1, 2\}\\ 0 \le P \le P_{max} \end{cases}$$
(9)

where H is the resource reservation probability, F is the candidate resource filtering coefficient, C is the reselection counter step offset, and P is the transmission power. 1 - H determines the probability of triggering resource reselect when RC is equal to 0. When the channel is not crowded, the higher the H, the lower the probability of triggering reselect, and the smaller the probability of resource collision caused by reselect. Cdetermines the reduced value of RC for each resource transfer using a predetermined resource, which defaults to 1. Dynamic reselection counter step offset can help SPS algorithm cope with more complex and changeable channel conditions, and a larger C can help SPS quickly enter the reselection time to avoid continuous resource collisions. P determines the range of broadcast transmission, and the appropriate transmission power can not only save energy but also reduce the interference between transmitters.

IV. RL BASED RESOURCE SELECTION ALGORITHM

This section mainly introduces the modeling of multi-agent environments and the enhanced SPS algorithm based on the QMIX algorithm.

A. Modeling of Multi-Agent Environments

Each vehicle interacts with the vehicular network environment as an agent and takes actions based on state observations, aiming to solve the optimization problem (9) At each time t, the vehicle takes a decision a_t based on the observed environment state o_t . The environment will be transferred to s_{t+1} , and then the vehicle will receive a reward r based on the decision made at the previous time. In the multi-agent environment architecture we established, the state observation space O, action space A, state space S, and reward function R are defined as follows:

1) Observation: In the actual environment under C-V2X Mode 4, the information that the vehicle can stably observe is very limited, mainly including two aspects: vehicle driving information and spectrum resource sensing information. This paper uses the following tuple to describe the observation information of the *t*th time step after vehicle *i* performs the action a_{t-1} at the *t*th time step:

$$o_t^i = \{Speed_t^i, Vec_t^i, \mathbf{CS}_t^i, RC_t^i\}$$
(10)

Vehicle driving information includes two elements: Speed and Vec, which are vehicle speed and vehicle driving direction respectively. These two pieces of information are instrument information that can be stably obtained when driving in the actual environment. Spectrum resource sensing information \mathbf{CS}_{t}^{i} : includes five elements: $CS1_{t}^{i}$, $CS2_{t}^{t}$, $CS3_{t}^{t}$, $CS4_{t}^{t}$, $CS5_{t}^{t}$. These five elements are the five statistics of the RSSI of the spectrum resources by the vehicle in the past broadcast period. After excluding the interference of half-duplex, these five statistics respectively represent the total number of resources, the number of resources with RSSI greater than the threshold, and the resource RSSI Cumulative sum, resource RSSI mean, and resource RSSI standard deviation. It is used to help the model perceive the overall situation of spectrum resources. In addition to the above observations of the environment, there is also a description of the protocol status. This paper uses the resource reselection counter RC_t^i to help the model understand the current status of the SPS protocol.

2) Action: :The action performed by vehicle i at time step t is defined as follows:

$$a_t^i = (H_t^i, F_t^i, C_t^i, P_t^i)$$
(11)

where H_t^i represents the dynamic reservation probability of spectrum resources. When the resource reservation counter returns to 0, whether to reserve resources will be decided based on the resource reservation probability. This paper uses this parameter to control the stability of resource reservations. F_t^i represents the proportion of candidate resources. Before resource reselection, it is necessary to filter the candidate pool according to the proportion of F_t^i , and then randomly select tile resources from the set. C_t^i is the offset step size of the resource reservation counter. In the SPS protocol, every time the vehicle performs a periodic broadcast, the resource reservation counter will be decremented by 1. When it returns



Fig. 1: QMIX-SPS resource allocation algorithm framework

to 0, resource reselection is possible. We use C_t^i to adjust the single change step size of the resource reservation technology. That is, each time a periodic broadcast is performed, the resource reservation counter is decremented by $(1 + C_t^i) \cdot P_t^i$ represents signal transmission power.

3) State: Under the centralized training and distributed execution architecture, global state information is only used during the training phase to help each agent model modify its strategy. Different from the state observation space, the training phase needs to use as comprehensive state information as possible to assist the agent in training. This information includes the following aspects: observation information of each vehicle, vehicle action information, and supplementary related information. Global status information at time t:

$$s_t = \{\mathbf{O}_t, \mathbf{a_{t-1}}, Location_t\}$$
(12)

where O_t is the collection of environmental observation data of all vehicles following the execution of action a_{t-1} , and a_{t-1} is the collection of actions that all vehicles select to do at the t-1th time step. The position of every vehicle at the tth time step, or *Location*_t, provides the extra relationship information. During the training phase, this parameter is intended to help the model measure the vehicle's distribution status.

4) Reward: The evaluation of the scheduling strategy in this paper is based on the SINR achievement rate PAS. The higher the PAS, the higher the proportion of communications that meet QoS constraints. r_t represents the global reward value for the t-1 event step as follows:

$$r_t = PAS_t - b_s \tag{13}$$

where b_s is the artificially set baseline. Since *PAS* is always positive, if PAS is used directly as a reward, the reward will always be positive, that is, any action selected will have a positive reward, which may cause difficulties in model training. Therefore, the *PAS* baseline b_s is introduced, and its value is the mean *PAS* value obtained by using the original SPS protocol under the same conditions. Essentially, it involves a type of comparative learning.

B. Dynamic Adaptive Parameters Optimization Algorithm QMIX-SPS

Based on multi-agent environment modeling, this paper proposes a joint dynamic parameters adaptive optimization algorithm for SPS based on the QMIX algorithm. The QMIX algorithm helps agents collaborate better to achieve global optimal results by constructing a hybrid network combined with the value function of a single agent. The overall architecture of the proposed algorithm is shown in Fig. 1. The proposed framework comprises two components: the agent network and the mixing network.

The agent network employs a deep recurrent Q network model to directly determine actions, integrating the recurrent unit (GRU) model within the deep recurrent neural network. This integration enables agents to leverage past trajectory data for decision-making. The agent network comprises three layers, with the GRU network situated in the middle, and the input and output layers consisting of fully connected neural networks. The input layer receives information such as the current observation data o_t^i of the vehicle and the previous action a_{t-1}^i taken. The GRU network requires the past hidden state output h_{t-1}^i of the agent as input and outputs the current hidden state h_t^i . In the output layer, the model generates the state-action value $Q(\tau_t^i, \cdot)$ for all actions based on the historical information τ_t^i , and employs a greedy strategy for action selection, as depicted in the subsequent equation:

$$\pi(\tau) = \arg\max_{a} Q_{\pi}(o, a) = a_t \tag{14}$$

The final value output by the agent network is expressed as $Q^i(\tau_t^i, a_t^i)$.

The mixing network is employed for value mixing, linking the global joint action value $Q_{tot}(\mathbf{o}, \mathbf{a})$ with the action value of each agent $Q^i(\tau_t^i, a_t^i)$. The mixing network, a basic two-layer feedforward neural network, is utilized to combine the action values of individual agents in a monotonic manner to generate $Q_{tot}(\mathbf{o}, \mathbf{a})$. To adhere to the QMIX constraint of monotonicity, the weight parameter of the hybrid network is constrained to be non-negative [13].

More specifically, the weight parameters and biases of each layer of the mixing network are produced by a distinct hypernetwork. Each layer of weights in the mixing network corresponds to a hypernetwork that takes the global state s as input and outputs the parameters of the feedforward neural network. The output is in vector form, which is then reshaped into a matching matrix based on predefined rules.

Furthermore, to enhance model stability and mitigate overestimation effects, the QMIX network integrates concepts from the traditional Deep Q Learning (DQN) algorithm, including empirical buffer pools and dual Q-networks. The overall endto-end loss function for the QMIX network can be expressed by the following equation:

$$\begin{cases} \mathcal{L}(\theta) = \sum_{i=1}^{b} \left[\left(y_{tot}^{i} - Q_{tot}(\boldsymbol{\tau}, \mathbf{a}, s; \theta) \right)^{2} \right] \\ y_{tot} = r + \delta \max_{a'} Q_{tot}(\boldsymbol{\tau'}, \mathbf{a'}, s'; \theta^{-}) \end{cases}$$
(15)

where b is the batch size sampled from the experience replay buffer for each training, and y_{tot} and θ^- represent the values obtained from the target network in the DQN as well as the network parameters. The main training process is shown in Algorithm 1.

V. SIMULATION RESULTS AND ANALYSIS

In this section, we introduce the simulation tools used in the study and the relevant experimental parameters. Then, the training process analysis, performance evaluation, and compatibility analysis of the proposed reinforcement learning algorithm are performed respectively.

A. Simulation Setup

Our simulation experiments are based on the open-source Python simulator Simulators-for-SPS [8]. We made some modifications to the simulator to make it more compliant with the TR 36.885 [14], mainly changing the channel model to Winner+B1 required by 3GPP. At the same time, the simulation experiments refer to the 3GPP C-V2X simulation guide [15]. We build an urban scenario with a $1299m \times 750m$ Manhattan grid. The scenario is composed of 3×3 units, and each road is two-way and four-lane. Vehicles move smoothly in the urban grid. We considered different vehicle densities and vehicle kinematics in the simulation and generated real traffic trajectory data by the road traffic simulator Simulation of Urban MObility (SUMO). The parameters of the simulation experiments are shown in Table I below.

B. Evaluation Metrics

To evaluate the performance of different resource selection algorithms, we introduce the following metrics in addition to PAS.

Algorithm 1: Training Algorithm

4

5

6

7

8

9

10 11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

```
Input: learning rate \alpha, replay buffer D, step limit
             step_{max}, episode limit episode_{max}, parameter
            synchronization interval step_{inr}, batch-size,
            reward factor \delta
   Output: \theta, the parameters of mixing network, agent
               networks and hypernetwork
1 Initialise \theta
2 step = 0, \theta^- = \theta
3 while step < step_{max} do
        t = 0
        Get initial state s_0
        while s_t \neq terminal and t \leq episode_{max} do
            foreach i in Vehicles do
                 Get available actions A_t^i for vehicle i
                 \tau_t^i = \tau_t^i \cup \{(o_t^i, a_t^i)\}
                 \epsilon=epsilon-annealing(step)
                 a_t^i =
                     argmax Q(\tau^i_t, a^i_t) with probability 1 - \epsilon
                       a_t^i \in A_t^i
                     Randomly select action from A_t^i
                      with probability \epsilon
            end
            Get reward r_t and next state s_{t+1}
            D=D\cup\{(s_t, a_t, r_t, s_{t+1})\}
            t = t + 1, step = step + 1
        end
        if D > batch-size then
            train-batch b \leftarrow random batch of episodes from
              D
            foreach b in train-batch do
                 Calculate Q_{tot} using Mixing-network with
                   Hypernetwork(s; \theta))
                 Calculate target Q_{tot} using Mixing-network
                   with Hypernetwork(s'; \theta^{-}))
             end
             y_{tot} = r + \delta max Q_{tot}(\tau', a', s'; \theta^{-})
             \Delta Q_{tot} = y_{tot} - Q_{tot}
             \Delta \theta = \nabla_{\theta} (\Delta Q_{tot})^2
            \theta = \theta - \alpha \Delta \theta
        end
        Synchronize parameters every step_{inr}: \theta^- \leftarrow \theta
29 end
```

- Packet delivery rate (PDR): The ratio of the number of packets successfully received by the vehicle to the number of packets expected to be received by the vehicle.
- Transmission speed: The transmission rate can be calcu-• lated from (4). This metric can be used to measure the instantaneous capacity of the link.

TABLE I: S	imulation	parameters
------------	-----------	------------

Parameter	Value	
Vehicle speed limit	60 km/h	
Vehicle average initial speed	36 km/h	
Vehicle acceleration	[-4.5 m/s ² ,4.5 m/s ²]	
Carrier frequency	5.9 GHz	
Channel bandwidth	10 MHz	
Subchannels per subframe	5	
RBs per subchannel	10	
Modulation and coding scheme	MCS 4	
Path loss model	WINNER+B1	
Antenna gain	3 dB	
Antenna height	1.5 m	
Shadow fading standard and deviation	3 dB, 4 dB	
RSRP threshold	-128 dBm	
Resource retention period	100ms	
Message sending frequency	10HZ	
Packet size	190 Bytes	

C. Simulation Result



Fig. 2: Training performance evaluation for QMIX-SPS

1) Training process analysis: Fig. 2(a) shows the training performance of our proposed QMIX-SPS algorithm with a mini-batch size of 256. It is observed that the loss function value of the mixing network decreases rapidly with the increase of training episodes until it converges to a minimal value. Fig. 2(b) shows the changing trend of the average reward as the number of training sets increases. The average reward climbs oscillate over time before reaching a dynamic equilibrium, which is a normal phenomenon in reinforcement

learning. It can be seen that the training algorithm we proposed performs well at convergence.

2) Performance of the Algorithm: As shown in Fig. 3, we investigate the performance of the SPS algorithm under varying vehicle densities in an urban scenario. When the total vehicle count is less than or equal to 500, the SPS algorithm's performance across metrics is not significantly affected by the vehicle count. However, once the vehicle count exceeds 500, we observe a relative decline in the SPS algorithm's performance, aligning with the limited available spectrum resources. This suggests the SPS algorithm's adaptive capabilities in low-density scenarios, but limited performance improvements in high-density settings. To address this, we propose the QMIX-SPS (RL) algorithm, which builds on the SPS mechanism with adaptive control of the algorithm parameter and signal transmit power. Experiments confirm the QMIX-SPS algorithm significantly enhances key C-V2X mode 4 communication metrics in the urban Manhattan grid scenario. Our findings highlight the need for scalable communication strategies to ensure reliable V2X performance in high-density urban environments. The QMIX-SPS algorithm offers a promising approach to improving V2X communication performance in complex urban settings.

3) Compatibility analysis: Since C-V2X is an evolving technology, different versions of MAC layer protocols must coexist. This paper considers that the newly proposed algorithm needs to maintain a certain degree of compatibility with the existing SPS protocol. The compatibility here refers to reducing the impact of introducing new resource allocation algorithms into the SPS environment in the original environment. Specifically, we consider different proportions of QMIX-SPS vehicles deployed in the SPS environment. The deployment proportion is equal to the ratio of the number of vehicles deploying QMIX-SPS to the total number of vehicles. As can be seen from Fig. 4, overall the QMIX-SPS algorithm has good compatibility with SPS, and the deployment of QMIX-SPS will not affect the performance of the SPS vehicle network. On the contrary, deploying QMIX-SPS can improve the performance of the SPS network. This performance improvement is particularly obvious in the case of low-density (100 vehicles) and high-density (600 vehicles), which illustrates the coordinating role of QMIX-SPS in vehicle communication networks. The performance optimization of the SPS network by QMIX-SPS is not comprehensive. For example, QMIX-SPS decreases the performance of the SPS network in an 80% deployment ratio setting with 400 vehicles.

VI. CONCLUSION

In this paper, we first introduce the spectrum resource allocation problem in C-V2X Mode 4 in detail and describe the system modeling and formal definition of this problem. Before formally introducing the algorithm proposed in this paper, we explain the design of reinforcement learning elements for this problem, including the design of action space, state space, observation space, and reward function. Subsequently, this



Fig. 3: Performance of QMIX-SPS.



Fig. 4: Compatibility of QMIX-SPS.

section introduces the proposed QMIX-SPS algorithm including the implementation process, optimization mechanism, and pseudo-code description. Finally, this section uses two groups of experiments to verify the impact of the algorithm on C-V2X Mode 4 communication performance and its compatibility with the original SPS algorithm. Experimental results prove that our proposed algorithm can effectively improve global communication performance while maintaining good compatibility with the SPS. In future work, we will further analyze the challenges of parameter optimization using reinforcement learning.

REFERENCES

- [1] A. R. Khan, M. F. Jamlos, N. Osman, M. I. Ishak, F. Dzaharudin, Y. K. Yeow, and K. A. Khairi, "Dsrc technology in vehicle-to-vehicle (v2v) and vehicle-to-infrastructure (v2i) iot system for intelligent transportation system (its): A review," *Recent Trends in Mechatronics Towards Industry 4.0: Selected Articles from iM3F 2020, Malaysia*, pp. 97–106, 2022.
- [2] 3GPP, "Evolved Universal Terrestrial Radio Access (E-UTRA); Radio Resource Control (RRC); Protocol specification," 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 36.300, 12 2016, version 14.1.0.
- [3] M. Gonzalez-Martín, M. Sepulcre, R. Molina-Masegosa, and J. Gozalvez, "Analytical models of the performance of c-v2x mode 4 vehicular communications," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 2, pp. 1155–1166, 2018.
- [4] Z. Ali, S. Lagén, L. Giupponi, and R. Rouil, "3gpp nr v2x mode 2: Overview, models and system-level evaluation," *IEEE Access*, vol. 9, pp. 89554–89579, 2021.
- [5] X. Gu, J. Peng, L. Cai, X. Zhang, and Z. Huang, "Markov analysis of c-v2x resource reservation for vehicle platooning," in 2022 IEEE 95th Vehicular Technology Conference: (VTC2022-Spring). IEEE, 2022, pp. 1–5.

- [6] V. Todisco, S. Bartoletti, C. Campolo, A. Molinaro, A. O. Berthet, and A. Bazzi, "Performance analysis of sidelink 5g-v2x mode 2 through an open-source simulator," *IEEE Access*, vol. 9, pp. 145648–145661, 2021.
- [7] A. Dayal, V. K. Shah, B. Choudhury, V. Marojevic, C. Dietrich, and J. H. Reed, "Adaptive semi-persistent scheduling for enhanced on-road safety in decentralized v2x networks," in 2021 *IFIP Networking Conference* (*IFIP Networking*). IEEE, 2021, pp. 1–9.
- [8] X. Gu, J. Peng, L. Cai, Y. Cheng, X. Zhang, W. Liu, and Z. Huang, "Performance analysis and optimization for semi-persistent scheduling in c-v2x," *IEEE Transactions on Vehicular Technology*, vol. 72, no. 4, pp. 4628–4642, 2022.
- [9] T. Kim, Y. Kim, M. Jung, and H. Son, "Intelligent partial sensing based autonomous resource allocation for nr v2x," *IEEE Internet of Things Journal*, 2023.
- [10] B. Wang, J. Zheng, N. Mitton, and C. Li, "Inp-crs: an intra-platoon cooperative resource selection scheme for c-v2x networks," *IEEE Communications Letters*, 2023.
- [11] L. Liang, H. Ye, and G. Y. Li, "Spectrum sharing in vehicular networks based on multi-agent reinforcement learning," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 10, pp. 2282–2292, 2019.
- [12] M. Parvini, M. R. Javan, N. Mokari, B. Abbasi, and E. A. Jorswieck, "Aoi-aware resource allocation for platoon-based c-v2x networks via multi-agent multi-task reinforcement learning," *IEEE Transactions on Vehicular Technology*, 2023.
- [13] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson, "Monotonic value function factorisation for deep multiagent reinforcement learning," *Journal of Machine Learning Research*, vol. 21, no. 178, pp. 1–51, 2020.
- [14] 3GPP, "Study on Ite-based v2x services," 3rd Generation Partnership Project (3GPP), Technical Specification (TS) 36.885, 07 2016, version 14.0.0.
- [15] M. H. C. Garcia, A. Molina-Galan, M. Boban, J. Gozalvez, B. Coll-Perales, T. Şahin, and A. Kousaridas, "A tutorial on 5g nr v2x communications," *IEEE Communications Surveys & Tutorials*, vol. 23, no. 3, pp. 1972–2026, 2021.