

Multi-Timescale Control and Communications with Deep Reinforcement Learning—Part II: Control-Aware Radio Resource Allocation

Lei Lei, *Senior Member, IEEE*, Tong Liu, Kan Zheng, *Senior Member, IEEE*, Xuemin (Sherman) Shen *Fellow, IEEE*

Abstract— In Part I of this two-part paper (Multi-Timescale Control and Communications with Deep Reinforcement Learning—Part I: Communication-Aware Vehicle Control), we decomposed the multi-timescale control and communications (MTCC) problem in Cellular Vehicle-to-Everything (C-V2X) system into a communication-aware Deep Reinforcement Learning (DRL)-based platoon control (PC) sub-problem and a control-aware DRL-based radio resource allocation (RRA) sub-problem. We focused on the PC sub-problem and proposed the MTCC-PC algorithm to learn an optimal PC policy given an RRA policy. In this paper (Part II), we first focus on the RRA sub-problem in MTCC assuming a PC policy is given, and propose the MTCC-RRA algorithm to learn the RRA policy. Specifically, we incorporate the PC advantage function in the RRA reward function, which quantifies the amount of PC performance degradation caused by observation delay. Moreover, we augment the state space of RRA with PC action history for a more well-informed RRA policy. In addition, we utilize reward shaping and reward backpropagation prioritized experience replay (RBPER) techniques to efficiently tackle the multi-agent and sparse reward problems, respectively. Finally, a sample- and computational-efficient training approach is proposed to jointly learn the PC and RRA policies in an iterative process. In order to verify the effectiveness of the proposed MTCC algorithm, we performed experiments using real driving data for the leading vehicle, where the performance of MTCC is compared with those of the baseline DRL algorithms.

Index Terms—Multi-Timescale Decision-Making; Radio Resource Allocation; Deep Reinforcement Learning

I. INTRODUCTION

In Part I of this two-part paper, we introduce the problem of multi-timescale control and communications (MTCC) in Cellular Vehicle-to-Everything (C-V2X) systems, where two types of decisions need to be made at different timescales: (1) platoon control (PC) decisions that are made with a coarse time grid of every T ms control interval; (2) radio resource allocation (RRA) decisions that are made with a fine time grid of every 1 ms communication interval. PC aims to determine the control inputs for following vehicles so that all vehicles move at the same speed while maintaining the desired distance between each pair of preceding vehicle (i.e., predecessor) and following vehicle (i.e., follower) [1]. The C-V2X communications enable information exchange between vehicles so that more informed PC decisions can be made to reduce the inter-vehicle distance while guaranteeing string stability. In the C-V2X system, Vehicle-to-Vehicle (V2V) links coexist with Vehicle-to-Infrastructure (V2I) links. A V2I link

connects a vehicle to the Base Station (BS) and is used for high-throughput services, and a V2V link connects a pair of predecessor and follower for periodic transmission of Collaborative Adaptive Message (CAM) according to the Predecessors Following (PF) information typology (IFT) [2]. The PC decisions are made at each follower with the target of optimizing its PC performance, which is affected by the observation delay that depends on the RRA decisions of the C-V2X system. The RRA decisions including sub-channel allocation and power control are made at each predecessor with the targets of (1) maximizing the V2I throughput; (2) minimizing the PC performance degradation due to delayed observation. The two targets are contradictory with each other and an optimal trade-off should be struck. The trade-off heavily depends on the impact of observation delay on PC performance, which in turn is affected by the PC decisions.

The interplay between RRA and PC necessitates the collaborative design of communications and control functions. RRA in C-V2X systems should be control-aware, taking into account the control performance degradation due to the delay or packet loss for control-related information delivery. Meanwhile, PC should be communication-aware, considering the statistical properties of random delay and packet loss in C-V2X communications. Both PC and RRA are Sequential Stochastic Decision Problem (SSDP), where a sequence of decisions have to be made over a specific time horizon for a dynamic system whose states evolve in the face of uncertainty. We believe that tackling the SSDP under a unified Deep Reinforcement Learning (DRL) framework will better reveal the inter-dependency between control and communications, and thus facilitate the joint optimization task [3].

Since employing the full-space approach to solve the multi-timescale SSDP yields formidable computation complexity, we decompose the MTCC problem into two sub-problems, i.e., (1) communication-aware DRL-based PC, and (2) control-aware DRL-based RRA. In Part I of this two-part paper, we have studied the communication-aware DRL-based PC assuming an RRA policy is given, and proposed the MTCC-PC algorithm to learn the PC policy. The PC problem with observation delay is essentially a Random Delay Decentralized Partially Observable Markov Decision Process (RD-Dec-POMDP). Each follower is a PC agent, which makes a local and delayed observation at each control interval and decides on its local actions to maximize its expected cumulative individual reward. We augment the PC state space with the

observation delay and PC action history, and define the reward function for the augmented state to construct an augmented state Markov Decision Process (MDP). It is proved that the optimal policy for the MDP is optimal for the PC problem with delayed observation. We show by experiments that the proposed MTCC-PC algorithm outperforms the state-of-the-art communication-aware control by training in a delayed environment generated by fine-grained embedded simulation of C-V2X communications rather than by a simple stochastic delay model.

This paper is Part II of the two-part paper, in which we first focus on the control-aware DRL-based RRA problem assuming a PC policy is available, and propose the MTCC-RRA algorithm to learn the RRA policy. Different from most existing works on RRA in C-V2X systems [4]–[14], the RRA decisions are optimized taking into account the impact of observation delay on PC performance. The main contributions of this paper (Part II) are summarized as follows:

- **DRL-based Value of Information (VoI) Estimation and Optimization:** A systematic method is proposed for the DRL-based PC agent of each follower to estimate the VoI of any given information that are to be shared by its predecessor, where the VoI is leveraged by the DRL-based RRA agents of the predecessors to make optimal RRA decisions. Firstly, the PC performance degradation due to observation delay for each follower is quantified by the performance difference between the de-facto PC policy with observation delay and the optimal policy with undelayed observation. Then, the optimization objective of RRA is formulated by incorporating the performance degradation term. To solve such an optimization problem, it is rigorously proved that the advantage function of the optimal PC policy should be included in the RRA reward function. Moreover, the PC action history has to be included in the RRA state space to make the RRA policy more well-informed. One important advantage of this design is that the RRA decisions are made based on “VoI per control interval”, i.e., the advantage function of the optimal PC policy, which provides a finer-grained VoI compared with the existing VoI calculation methods.
- **Efficient DRL solution addressing multi-agent and sparse reward problems:** Since the RRA agents cooperatively optimize the global expected cumulative reward, and it is hard for each agent to deduce its own contribution to the global reward, the simple independent learner (IL) approach used in PC is not suitable for the RRA problem falling into the multi-agent domain. To solve this problem, we apply the reward shaping technique to design an individual RRA reward for each RRA agent. In addition, since the RRA rewards that reflect the PC performance degradation are only provided at the last communication interval of each control interval, it is difficult to learn an efficient RRA policy due to the sparse reward problem coupled with long-horizon planning. To solve this problem, we use the reward backpropagation prioritized experience replay (RBPER) technique [15] to improve training efficiency.

- **Sample- and computational-efficient training approach:** Since it is impossible to learn the optimal PC policy or RRA policy first without knowing the other optimal policy due to the interactions between PC and RRA, we design a sample- and computational-efficient training approach to jointly train MTCC-PC and MTCC-RRA in an iterative process to learn the optimal PC and RRA policies.

The rest of this paper is organized as follows. We recall the system model for the MTCC problem and the MTCC-PC algorithm in Section II. Subsequently, Sections III and IV introduce the control-aware DRL model for RRA and the MTCC-RRA algorithm. Then, a joint optimization algorithm for MTCC is presented in Section V. Section VI conducts the performance experiments, demonstrating that our proposed algorithm is superior to the baseline algorithms. Finally, Section VII concludes the paper.

II. SYSTEM MODEL AND PREVIOUS RESULTS

A. Multi-timescale Decision-Making Framework

Consider a platoon with a number of $N > 2$ vehicles, i.e., $\mathcal{V} = \{0, 1, \dots, N - 1\}$. All vehicles communicate with one another using C-V2X communications. The PC problem is considered within a finite time horizon, which is discretized into K equal-length control intervals indexed by $k \in \mathcal{K} = \{0, 1, \dots, K - 1\}$. The duration of each control interval is T milliseconds (ms). At each control interval k , the PC module of each follower $i \in \mathcal{V} \setminus \{0\}$ determines the vehicle control input $a_{i,k}^{\text{CL}}$ based on the observations of the system state. The vehicle driving status is sampled at time kT . We adopt the PF IFT, where the CAM $c_{i-1,k} = \{p_{i-1,k}, v_{i-1,k}, acc_{i-1,k}\}$ of the predecessor $i - 1 \in \mathcal{V} \setminus \{N - 1\}$ are transmitted to the follower i . For this purpose, each control interval k is further divided into T communication intervals of every 1 ms indexed by $t \in \mathcal{T} = \{0, 1, \dots, T - 1\}$ on a faster timescale. Dynamic scheduling is considered, where the C-V2X communication module makes RRA decisions to transmit CAM at each communication interval (k, t) .

Since the PC decisions are made with a coarse time grid of every T ms, while the RRA decisions are made with a fine time grid of every 1 ms, we have a multi-timescale decision-making problem.

B. Platoon Control Module and C-V2X Communications Module

For the platoon control module, each vehicle $i \in \mathcal{V}$ obeys the dynamics model described by a first-order system, and the dynamics model is derived in discrete time on the basis of forward Euler discretization.

$$p_{i,k+1} = p_{i,k} + Tv_{i,k}, \quad (1)$$

$$v_{i,k+1} = v_{i,k} + Tacc_{i,k}, \quad (2)$$

$$acc_{i,k+1} = \left(1 - \frac{T}{\tau_i}\right)acc_{i,k} + \frac{T}{\tau_i}a_{i,k}^{\text{CL}}, \quad (3)$$

where τ_i is a time constant representing driveline dynamics. The tracking errors, i.e., gap-keeping error $e_{pi,k}$ and velocity error $e_{vi,k}$ of follower i are defined as

$$e_{pi,k} = d_{i,k} - d_{r,i,k}, \quad (4)$$

$$e_{vi,k} = v_{i-1,k} - v_{i,k}, \quad (5)$$

where $d_{i,k} = p_{i-1,k} - p_{i,k} - L_{i-1}$ is the headway of follower i at control interval k , and $d_{r,i,k} = r_i + h_i v_{i,k}$ is the desired headway.

For C-V2X communications module, we consider a typical urban C-V2X network, where M V2I links (uplink considered) coexist with $N - 1$ V2V links. A V2I link $m \in \mathcal{M} = \{0, \dots, M - 1\}$ connects a vehicle to the Base Station (BS) and is used for high-throughput services. According to the PF IFT, a V2V link $i \in \mathcal{V} \setminus \{N - 1\}$ connects a pair of predecessor i and follower $i + 1$ for periodic transmission of CAM. In order to enhance spectrum utilization, the V2V links reuse the sub-channels of the V2I links for CAM transmission. We use the binary allocation indicator $\theta_{i,m,(k,t)} \in \{0, 1\}$ to indicate whether V2V link i occupies sub-channel m at communication interval (k, t) or not. In each communication interval (k, t) , each vehicle i transmits the data in its queue according to the local RRA decisions. Assume the instantaneous channel gain of V2V link i over sub-channel m (occupied by V2I link m) at communication interval (k, t) is denoted by $G_{i,m,(k,t)}$. Similarly, let $G_{m,(k,t)}$ denote the channel gain of the V2I link m ; $G_{i,B,m,(k,t)}$ the interference channel gain from V2V link i transmitter to V2I link m receiver; $G_{B,i,m,(k,t)}$ the interference channel gain from V2I link m transmitter to V2V link i receiver; and $G_{j,i,m,(k,t)}$ the interference channel gain from the V2V link j transmitter to the V2V link i receiver over the sub-channel m . The transmit power of V2V link i is $P_{i,m,(k,t)}^V$ and the constant transmit power of V2I link m is P_m^I . The SINR $\gamma_{m,(k,t)}$ of V2I link m and the SINR $\gamma_{i,m,(k,t)}$ of V2V link i on sub-channel m at communication interval (k, t) are derived by

$$\gamma_{m,(k,t)} = \frac{P_m^I G_{m,(k,t)}}{\sigma^2 + \sum_{i \in \mathcal{V} \setminus \{N-1\}} \theta_{i,m,(k,t)} P_{i,m,(k,t)}^V G_{i,B,m,(k,t)}}, \quad (6)$$

and

$$\gamma_{i,m,(k,t)} = \frac{P_{i,m,(k,t)}^V G_{i,m,(k,t)}}{\sigma^2 + I_{i,m,(k,t)}}, \quad (7)$$

respectively, where σ^2 is the power of channel noise which satisfies the independent Gaussian distribution with a zero mean value. $I_{i,m,(k,t)}$ is the total interference power received by V2V link i over sub-channel m , where

$$I_{i,m,(k,t)} = P_m^I G_{B,i,m,(k,t)} + \sum_{j \in \mathcal{V} \setminus \{i, N-1\}} \theta_{j,m,(k,t)} P_{j,m,(k,t)}^V G_{j,i,m,(k,t)}.$$

Then, the instantaneous data rate in terms of CAM of V2V link i and the instantaneous data rate of V2I link m at communication interval (k, t) are respectively derived as

$$r_{i,(k,t)}^{\text{CAM}} = \frac{\sum_{m=0}^{M-1} W \log_2(1 + \gamma_{i,m,(k,t)})}{N_c}, \quad (8)$$

and

$$r_{m,(k,t)} = W \log_2(1 + \gamma_{m,(k,t)}), \quad (9)$$

where W is the bandwidth of a sub-channel, and N_c is the constant CAM size.

Let $q_{i,(k,t)}^{\text{CAM}}$ denote the CAM queue length of vehicle i at communication interval (k, t) . If the queue length $q_{i,(k,t)}^{\text{CAM}}$ reaches the buffer capacity N_Q , the subsequent arriving data will be dropped. The queue process evolves as

$$q_{i,(k,t+1)}^{\text{CAM}} = \begin{cases} \min \left[N_Q, \max[0, q_{i,(k,t)}^{\text{CAM}} - 10^{-3} \times r_{i,(k,t)}^{\text{CAM}}] + 1 \right], & \text{if } t = 0 \\ \max \left[0, q_{i,(k,t)}^{\text{CAM}} - 10^{-3} \times r_{i,(k,t)}^{\text{CAM}} \right], & \text{otherwise} \end{cases} \quad (10)$$

The CAM that is not transmitted during control interval k will continue to be transmitted in the next control interval $k + 1$.

C. Correlation between Platoon Control Decisions and Radio Resource Allocation Decisions

Let $\tau_{i,k}$ be the observation delay of follower i at control interval k . Thus, $c_{i-1,k-\tau_{i,k}} = \{p_{i-1,k-\tau_{i,k}}, v_{i-1,k-\tau_{i,k}}, acc_{i-1,k-\tau_{i,k}}\}$ is the most recent available delayed CAM at follower i , which correspond to the position, velocity, and acceleration sampled at predecessor $i - 1$ in control interval $k - \tau_{i,k}$. Therefore, the observed driving status of vehicle i is defined as

$$x_{i,k-\tau_{i,k}} = \{e_{pi,k-\tau_{i,k}}, e_{vi,k-\tau_{i,k}}, acc_{i,k-\tau_{i,k}}, acc_{i-1,k-\tau_{i,k}}\}, \quad (11)$$

where $e_{pi,k-\tau_{i,k}} = p_{i-1,k-\tau_{i,k}} - p_{i,k-\tau_{i,k}} - L_{i-1} - d_{r,i,k}$ and $e_{vi,k-\tau_{i,k}} = v_{i-1,k-\tau_{i,k}} - v_{i,k-\tau_{i,k}}$.

The observation delay $\tau_{i,k}$ depends on the transmission delay of CAM over V2V link $i - 1$, which can be derived from $q_{i-1,(k-1,T)}^{\text{CAM}}$ or $q_{i-1,(k,0)}^{\text{CAM}}$ as

$$\tau_{i,k} = \lceil q_{i-1,(k-1,T)}^{\text{CAM}} \rceil + 1 = \lceil q_{i-1,(k,0)}^{\text{CAM}} \rceil + 1. \quad (12)$$

A detailed analysis of the correlation between PC and RRA decisions can be found in Part I of this two-part paper.

D. Communication-Aware PC and the Corresponding DRL Solution

We assume that the RRA policy π^{CM} is available and focus on learning the PC policy π_i^{CL} , $i \in \mathcal{V} \setminus \{0\}$.

1) Communication-aware PC model:

a) *PC state*: The state for each PC agent i at control interval k is defined as

$$S_{i,k}^{\text{CL}} = \{x_{i,k-\tau_{i,k}}, \{a_{i,k'}^{\text{CL}}\}_{k'=k-\tau_{\max}}, \tau_{i,k}\}. \quad (13)$$

b) *PC action*: The control input, $a_{i,k}^{\text{CL}}$ of the PC agent i is regarded as its PC action at control interval k .

c) *PC reward function*: The individual reward for each PC agent i is given by

$$R_{i,k}^{\text{CL}}(x_{i,k}, a_{i,k}^{\text{CL}}) = -\left\{ \left| \frac{e_{p,i,k}}{\hat{e}_{p,\max}} \right| + \alpha_1 \left| \frac{e_{v,i,k}}{\hat{e}_{v,\max}} \right| + \alpha_2 \left| \frac{a_{i,k}^{\text{CL}}}{a_{\max}^{\text{CL}}} \right| + \alpha_3 \left| \frac{j_{i,k}}{2\text{acc}_{\max}/T} \right| \right\}, \quad (14)$$

where $j_{i,k} = \frac{\text{acc}_{i,k+1} - \text{acc}_{i,k}}{T} = -\frac{1}{\tau_i} \text{acc}_{i,k} + \frac{1}{\tau_i} a_{i,k}^{\text{CL}}$ is the jerk.

The objective of the PC problem is for each PC agent i to find the optimal policy $\pi_i^{\text{CL}*}$ under delayed observation that maximizes its individual expected return J_i^{CL} , i.e.,

$$\pi_i^{\text{CL}*} = \arg \max_{\pi_i^{\text{CL}}} J_i^{\text{CL}}, \quad \forall i \in \mathcal{V} \setminus \{0\}. \quad (15)$$

where

$$J_i^{\text{CL}} = \mathbb{E}_{\pi_{\text{CM}}} \mathbb{E}_{\pi_i^{\text{CL}}} \left[\sum_{k=0}^{K-1} \gamma^k R_{i,k}^{\text{CL}} \right], \quad 0 \leq \gamma \leq 1. \quad (16)$$

Note that γ is the PC reward discount factor.

2) *DRL solution for communication-aware PC model*:

a) *Multi-agent problem in DRL-based PC*: The PC problem corresponds to a Dec-POMDP and lies in the multi-agent domain. We adopt the IL approach where each agent learns independently since the non-stationary environment issue for IL is trivial in the PC problem and the credit assignment issue in the multi-agent problem does not exist for our PC model.

b) *Random observation delay problem in DRL-based PC*: We approximately consider that the undelayed driving status $x_{i,k}$ at PC agent i is Markov, ignoring the impact of the predecessors' actions on $x_{i,k+1}$. However, each PC agent i can only observe the delayed driving status $x_{i,k-\tau_{i,k}}$ instead of $x_{i,k}$. It is proved in Theorem 1 of Part I that $S_{i,k}^{\text{CL}}$ becomes a Markov state by augmenting the delayed observation of driving status with action history.

We construct an MDP $\tilde{\mathcal{M}}_i = (S_{i,k}^{\text{CL}}, a_{i,k}^{\text{CL}}, \tilde{R}_{i,k}^{\text{CL}}, p, \gamma)$ for the delayed observations, where the reward function $\tilde{R}_{i,k}^{\text{CL}}(S_{i,k}^{\text{CL}}, a_{i,k}^{\text{CL}})$ for the augmented state of each follower $i \in \mathcal{V} \setminus \{0\}$ is defined as $\tilde{R}_{i,k}^{\text{CL}}(S_{i,k}^{\text{CL}}, a_{i,k}^{\text{CL}}) = \mathbb{E}_{x_{i,k}} [R_{i,k}^{\text{CL}}(x_{i,k}, a_{i,k}^{\text{CL}}) | S_{i,k}^{\text{CL}}]$. It is proved in Theorem 2 of Part I that the optimal policy $\tilde{\pi}_i^{\text{CL}*}$ for the augmented state MDP $\tilde{\mathcal{M}}_i$ is the same as the optimal policy $\pi_i^{\text{CL}*}$ in (15) for our PC problem under delayed observation.

c) *MTCC-PC Algorithm*: Based on the above discussion, the DDPG algorithm [16] is utilized to solve the PC problem. The Deterministic Policy Gradient (DPG) Theorem can be directly applied since $\tilde{\mathcal{M}}_i$ is an MDP, i.e., the deterministic policy gradient for the augmented state MDP $\tilde{\mathcal{M}}_i$ is

$$\nabla_{\theta_i^{\mu}} J_i^{\text{CL}}(\mu_i^{\text{CL}}) = \mathbb{E} [\nabla_{\theta_i^{\mu}} \mu_i^{\text{CL}}(S_{i,k}^{\text{CL}} | \theta_i^{\mu}) \nabla_a Q_i^{\text{CL}}(S_{i,k}^{\text{CL}}, a | \theta_i^Q) |_{a=\mu_i^{\text{CL}}(S_{i,k}^{\text{CL}} | \theta_i^{\mu})}]. \quad (17)$$

In order to sample the deterministic policy gradient in (17), we need to evaluate the action-value function $Q_{\mu_{\theta_i}}^{\text{CL}}(S_{i,k}^{\text{CL}}, a_{i,k}^{\text{CL}})$ of the augmented state MDP $\tilde{\mathcal{M}}_i$. Based on the Bellman equation, the expected cumulative discounted reward from control interval k is defined as

$$Q_{\mu_{\theta_i}}^{\text{CL}}(S_{i,k}^{\text{CL}}, a_{i,k}^{\text{CL}}) = \mathbb{E}_{x_{i,k}} [R_{i,k}^{\text{CL}}(x_{i,k}, a_{i,k}^{\text{CL}}) | S_{i,k}^{\text{CL}}]$$

$$+ \gamma \mathbb{E}_{S_{i,k+1}^{\text{CL}}} [Q_{\mu_{\theta_i}}^{\text{CL}}(S_{i,k+1}^{\text{CL}}, \mu_{\theta_i}(S_{i,k+1}^{\text{CL}})) | S_{i,k}^{\text{CL}}, a_{i,k}^{\text{CL}}]. \quad (18)$$

During training, the sampled deterministic policy gradient ascent on $Q_i^{\text{CL}}(S_{i,k}, \mu_i^{\text{CL}}(S_{i,k} | \theta_i^{\mu})) | \theta_i^Q$ with regard to θ_i^{μ} is used to train the actor network, and the critic network is trained by minimizing the Root Mean Square Error (RMSE) $L_{i,k} = y_{i,k} - Q_i^{\text{CL}}(S_{i,k}^{\text{CL}}, a_{i,k}^{\text{CL}} | \theta_i^Q)$ using the sampled gradient descent with respect to θ_i^Q , where the TD target $y_{i,k}$ is calculated by

$$y_{i,k} = R_{i,k}^{\text{CL}}(x_{i,k}, a_{i,k}^{\text{CL}}) + \gamma Q_{\mu_{\theta_i}}^{\text{CL}}(S_{i,k+1}^{\text{CL}}, \mu_{\theta_i}(S_{i,k+1}^{\text{CL}})). \quad (19)$$

Note that in the last control interval $K-1$ of the control episode, we still calculate $y_{i,k}$ using (19) instead of only the immediate reward. In this way, we convert the finite horizon SSDP to an infinite horizon one, since DDPG and most DRL algorithms are designed for infinite or indefinite horizon problems [17].

It is assumed that the PC agent i can sample the undelayed reward $R_{i,k}^{\text{CL}}(x_{i,k}, a_{i,k}^{\text{CL}})$ during training. This is possible since learning can be performed in a simulator or a laboratory in which the undelayed reward is available.

An important characteristic of the proposed MTCC-PC algorithm is that it is trained in a delayed environment generated by the simulation of C-V2X communications with de facto RRA policy rather than by a coarse-grained stochastic delay model. This is to ensure the delay distribution in the training environment is the same as that in the execution environment.

III. CONTROL-AWARE DRL-BASED RADIO RESOURCE ALLOCATION

We assume that the PC policy is available and focus on learning the RRA policy in this section.

The RRA problem is essentially a Dec-POMDP. Each predecessor $i \in \mathcal{V} \setminus \{N-1\}$ is an RRA agent, which makes a local observation at each communication interval (k, t) , and decides on its local actions to maximize the expected cumulative *global* reward. Since one of the objectives of RRA is to minimize the PC performance degradation due to delayed observation, the main challenges in formulating the DRL model are the design of the reward function and state space.

A. PC Performance Degradation Due to Delayed Observation

We quantify the PC performance degradation for each follower $i \in \mathcal{V} \setminus \{0\}$ by the performance difference between $\pi_i^{\text{CL}}(S_{i,k}^{\text{CL}})$ and $\hat{\pi}_i^{\text{CL}*}(x_{i,k})$, where the former is the policy in Section II when the PC agent only has access to the augmented state $S_{i,k}^{\text{CL}}$ with delayed observation; while the latter is the optimal policy when the PC agent has oracle access to the current driving status $x_{i,k}$ at each control interval k . Therefore, the total PC performance degradation ΔJ^{CL} over all the followers is

$$\Delta J^{\text{CL}} = \sum_{i=0}^{N-2} \left(J_{i+1}^{\text{CL}} - \hat{J}_{i+1}^{\text{CL}*} \right), \quad (20)$$

where J_{i+1}^{CL} is the expected return of policy $\pi_{i+1}^{\text{CL}}(S_{i+1,k}^{\text{CL}})$ given in (16), and

$$\hat{J}_{i+1}^{\text{CL}*} = \mathbb{E}_{\hat{\pi}_{i+1}^{\text{CL}*}} \left[\sum_{k=0}^{K-1} \gamma^k R_{i+1,k}^{\text{CL}} \right] \quad (21)$$

is the expected return of optimal policy $\hat{\pi}_i^{\text{CL}*}(x_{i+1,k})$, which is not affected by π^{CM} since it is assumed that the PC agent can observe the current driving status $x_{i+1,k}$ without delay.

B. RRA Reward Modeling

The optimization objective of RRA is formulated as

$$J^{\text{CM}} = \kappa_1 \sum_{m=0}^{M-1} \mathbb{E}_{\pi^{\text{CM}}} \left[\sum_{k=0}^{K-1} \sum_{t=0}^{T-1} \eta^{kT+t} r_{m,(k,t)} \right] + \kappa_2 \Delta J^{\text{CL}}, \quad (22)$$

where η is the RRA reward discount factor. The first term corresponds to the discounted sum throughput of all V2I links $m \in \mathcal{M}$ over all the communication intervals of a control episode, and the second term corresponds to the PC performance degradation given in (20). The weight factors κ_1 and κ_2 indicate the relative importance of minimizing PC performance degradation versus maximizing the V2I throughput.

The optimization objective J^{CM} should correspond to the expected return of the DRL model, i.e.,

$$J^{\text{CM}} = \mathbb{E}_{\pi^{\text{CM}}} \mathbb{E}_{\pi^{\text{CL}}} \left[\sum_{k=0}^{K-1} \sum_{t=0}^{T-1} \eta^{kT+t} R_{(k,t)}^{\text{CM}} \right], 0 \leq \eta \leq 1, \quad (23)$$

where the expectation is taken with respect to the probability distribution of the state-action trajectories when the RRA policy is π^{CM} and PC policy is π^{CL} . Since both RRA and PC are multi-agent problems, the RRA (resp. PC) policy corresponds to the set of policies of all the RRA (resp. PC) agents, i.e., $\pi^{\text{CM}} = \{\pi_i^{\text{CM}}\}_{i \in \mathcal{V} \setminus \{N-1\}}$ and $\pi^{\text{CL}} = \{\pi_i^{\text{CL}}\}_{i \in \mathcal{V} \setminus \{0\}}$.

We will try to derive the RRA reward function $R_{(k,t)}^{\text{CM}}$ from (22) and (23). Specifically, we have

$$R_{(k,t)}^{\text{CM}} = \begin{cases} \kappa_1 R_{\text{I},(k,t)}, & 0 \leq t < T-1 \\ \kappa_1 R_{\text{I},(k,t)} + \kappa_2 R_{\text{V},(k,T)}, & t = T-1 \end{cases}, \quad (24)$$

where

$$R_{\text{I},(k,t)} = \sum_{m=0}^{M-1} r_{m,(k,t)} \quad (25)$$

is the component related to V2I throughput, and $R_{\text{V},(k,T)}$ is the component related to the PC performance degradation.

At first glance, $R_{\text{V},(k,T)}$ seems to be the difference between the two PC rewards $R_{i+1,k+1}^{\text{CL}}(x_{i+1,k+1}, \pi_{i+1}^{\text{CL}}(S_{i+1,k+1}^{\text{CL}}))$ and $R_{i+1,k+1}^{\text{CL}}(x_{i+1,k+1}, \hat{\pi}_{i+1}^{\text{CL}*}(x_{i+1,k+1}))$ based on (20), (16), and (21). However, this is not correct since the expectations in (16) and (21) are with respect to different trajectory distributions. In order to derive $R_{\text{V},(k,T)}$, the following lemma is used.

Lemma 1: The PC performance degradation due to observation delay for follower i equals to the expected cumulative advantage of $\pi_i^{\text{CL}}(S_{i,k}^{\text{CL}})$ over $\hat{\pi}_i^{\text{CL}*}(x_{i,k})$, i.e.,

$$J_i^{\text{CL}} - \hat{J}_i^{\text{CL}*} = \mathbb{E}_{\pi^{\text{CM}}} \mathbb{E}_{\pi_i^{\text{CL}}} \left[\sum_{k=0}^{K-1} \gamma^k A_{\hat{\pi}_i^{\text{CL}*}}(x_{i,k}, \pi_i^{\text{CL}}(S_{i,k}^{\text{CL}})) \right], \quad (26)$$

$0 \leq \gamma \leq 1,$

where

$$\begin{aligned} & A_{\hat{\pi}_i^{\text{CL}*}}(x_{i,k}, \pi_i^{\text{CL}}(S_{i,k}^{\text{CL}})) \\ &= Q_{\hat{\pi}_i^{\text{CL}*}}(x_{i,k}, \pi_i^{\text{CL}}(S_{i,k}^{\text{CL}})) - V_{\hat{\pi}_i^{\text{CL}*}}(x_{i,k}) \end{aligned} \quad (27)$$

is the advantage function of policy $\hat{\pi}_i^{\text{CL}*}(x_{i,k})$.

The proof of Lemma 1 is given in Appendix A.

Based on Lemma 1, $R_{\text{V},(k,T)}$ is defined as

$$R_{\text{V},(k,T)} = \sum_{i=0}^{N-2} A_{\hat{\pi}_{i+1}^{\text{CL}*}}(x_{i+1,k+1}, \pi_{i+1}^{\text{CL}}(S_{i+1,k+1}^{\text{CL}})) \quad (28)$$

Since the advantage function $A_{\hat{\pi}_{i+1}^{\text{CL}*}}(x_{i+1,k}, \pi_{i+1}^{\text{CL}}(S_{i+1,k}^{\text{CL}}))$ belongs to the optimal policy $\hat{\pi}_{i+1}^{\text{CL}*}(x_{i+1,k})$ without observation delay, the PC agent $i+1$ needs to learn $\hat{\pi}_{i+1}^{\text{CL}*}(x_{i+1,k})$ and the corresponding advantage function in addition to the acting policy $\pi_{i+1}^{\text{CL}}(S_{i+1,k}^{\text{CL}})$.

Finally, the following theorem formally justifies the RRA reward function defined above.

Theorem 1: Given the reward function defined in (24), (25), and (28), the expected return as calculated in (23) is the same as the optimization objective in (22), where the relationship between the PC reward discount factor and RRA reward discount factor is $\eta = \gamma^{\frac{1}{T}}$.

The proof of Theorem 1 is given in Appendix B.

C. RRA State Definition

The state of each RRA agent i is defined as

$$S_{i,(k,t)}^{\text{CM}} = \{G_{i,(k,t)}, q_{i,(k,t)}^{\text{CAM}}, \{a_{i,k'}^{\text{CL}}\}_{k'=k-\tau_{\text{max}}+1}, t, \epsilon\}, \quad (29)$$

which includes four types of information as discussed below.

a) *Channel state information:* $G_{i,(k,t)}$ is the local observation of CSI at vehicle i in communication interval (k,t) , i.e.,

$$G_{i,(k,t)} = \{G_{i,m,(k,t)}, G_{j,i,m,(k,t)}, G_{B,i,m,(k,t)}, G_{i,B,m,(k,t)}, G_{m,(k,t)}\}_{m \in \mathcal{M}}. \quad (30)$$

Specifically, $G_{i,(k,t)}$ contains the following information: (i) the channel gain of V2V link i over all sub-channels $m \in \mathcal{M}$, $\{G_{i,m,(k,t)}\}_{m \in \mathcal{M}}$; (ii) the interference channel gain from all V2I link m , $\{G_{B,i,m,(k,t)}\}_{m \in \mathcal{M}}$, and other V2V link j , $\{G_{j,i,m,(k,t)}\}_{j \neq i, m \in \mathcal{M}}$; (iii) the interference channel gain from V2V link i to the BS over all sub-channel $m \in \mathcal{M}$, $\{G_{i,B,m,(k,t)}\}_{m \in \mathcal{M}}$; (iv) the channel gain of all V2I link m , $\{G_{m,(k,t)}\}_{m \in \mathcal{M}}$. Channel gains (i) and (ii) can be accurately estimated by the receiver of V2V link i at the beginning of each communication interval (k,t) [18], while (iii) and (iv) are estimated at the BS in each communication interval (k,t) and then broadcast to all vehicles in its coverage, incurring a small signaling overhead.

b) *Queue state information:* $S_{i,(k,t)}^{\text{CM}}$ includes the queue length $q_{i,(k,t)}^{\text{CAM}}$ at vehicle i , whose value at communication interval $(k,0) = (k-1,T)$ reflects the observation delay of vehicle $i+1$ at control interval k according to (12).

c) *PC information*: Since the RRA policy aims at minimizing PC performance degradation due to observation delay, it is our hypothesis that including some PC information in the RRA state could help the RRA agent make more informed decisions. However, it is not a trivial task to determine what information should be included. Ignoring useful information is obviously undesirable, while including useless information should also be avoided due to the curse-of-dimensionality issue.

Based on the following lemmas, Theorem 2 reveals the PC information that should be included in the RRA state.

Lemma 2: The optimal PC policy $\pi_i^{\text{CL}*}(S_{i,k}^{\text{CL}})$ based on delayed observation $S_{i,k}^{\text{CL}}$ is equivalent to the optimal PC policy $\hat{\pi}_i^{\text{CL}*}(x_{i,k})$ based on current driving status if $x_{i,k}$ can be fully determined by the delayed observation $S_{i,k}^{\text{CL}}$, i.e., $x_{i,k} = f(S_{i,k}^{\text{CL}})$ where $f(S_{i,k}^{\text{CL}})$ is a deterministic function of $S_{i,k}^{\text{CL}}$.

The proof of Lemma 2 is given in Appendix C.

Lemma 3: The current driving status $x_{i,k}$ can be fully determined by the delayed observation $S_{i,k}^{\text{CL}}$ if the sequence of actions of vehicle $i-1$ within the time window of $[k - \tau_{i,k}, k - 1]$ is available to vehicle i , i.e., $x_{i,k} = f(S_{i,k}^{\text{CL}}, \{a_{i-1,k'}^{\text{CL}}\}_{k'=k-\tau_{i,k}}^{k-1})$.

The proof of Lemma 3 is given in Appendix D.

Theorem 2: The impact of observation delay $\tau_{i,k}$ to the PC performance of vehicle i depends on the variation of vehicle $i-1$'s control action $a_{i-1,k'}^{\text{CL}}$ within the period $k' \in [k - \tau_{\max} + 1, k]$.

The proof of Theorem 2 is given in Appendix E.

According to Theorem 2, $\{a_{i,k'}^{\text{CL}}\}_{k'=k-\tau_{\max}+1}^k$ is included in $S_{i,(k,t)}^{\text{CM}}$ since it determines how sensitive the PC performance is to the observation delay for a specific control interval k . In other words, the characteristics of $a_{i-1,k'}^{\text{CL}}$ within the period $k' \in [k - \tau_{\max} + 1, k]$ determines the VoI of the received message at vehicle i at control interval k .

Remark 1 (VoI per control interval): In the context of NCS or task-oriented scheduling, Age of Information (AoI) and VoI are two important metrics. AoI captures the importance of information by measuring its timeliness attribute [19], while VoI measures how much the recipient of the information can reduce the uncertainty of the stochastic processes related to decision-making [20]–[22]. It is shown in [23] that VoI is related to AoI while providing more accurate guidance to the scheduler design. However, VoI is normally derived based on mutual information as a single value for the whole control task, ignoring the variation in the significance of transmitted data during the control episode. In this paper, we design the RRA reward function to capture VoI based on the PC performance degradation due to the observation delay of shared information in each control interval. Moreover, we include the PC information in the RRA state, so that the RRA decisions are made based on the *VoI per control interval* that captures VoI at a finer time granularity.

d) *Other information*: In order to stabilize experience replay for multi-agent DRL, we augment $S_{i,(k,t)}^{\text{CM}}$ with the exploration rate ϵ [24], i.e., the probability of random selecting action in ϵ -greedy policy for Deep Q Networks (DQN).

Finally, the index of communication interval t is a part of the RRA state since it relates to the reward function in (24).

D. RRA Action

The RRA action of each vehicle $i \in \mathcal{V} \setminus \{N-1\}$ at communication interval (k, t) includes the decisions on the sub-channel allocation $\{\theta_{i,m,(k,t)}\}_{m \in \mathcal{M}}$ and transmit power $\{P_{i,m,(k,t)}^{\text{V}}\}_{m \in \mathcal{M}}$ for V2V link i . Since we consider that at most one sub-channel is allocated to a V2V link at any communication interval (k, t) , i.e., $\sum_{m=0}^{M-1} \theta_{i,m,(k,t)} \leq 1$, we can simplify the sub-channel allocation action from $\{\theta_{i,m,(k,t)}\}_{m \in \mathcal{M}}$ to $m_i \in \mathcal{M} \cup \{-1\}$, where $m_i \in \mathcal{M}$ means that sub-channel m_i is selected and $m_i = -1$ means that no sub-channel is allocated to V2V link i .

We discretize the transmit power to four levels for ease of learning and practical circuit restriction [9], [11], i.e., $P_{i,m,(k,t)}^{\text{V}} \in \mathcal{A}_P = \{23, 15, 5, -100\}$ dBm. Noted that -100 dBm can be considered as zero transmit power.

Therefore, the RRA action of each vehicle $i \in \mathcal{V} \setminus \{N-1\}$ at communication interval (k, t) is defined as

$$a_{i,(k,t)}^{\text{CM}} = \left\{ m_i, P_{i,m,(k,t)}^{\text{V}} \mid m_i \in \mathcal{M} \cup \{-1\}, P_{i,m,(k,t)}^{\text{V}} \in \mathcal{A}_P \right\}. \quad (31)$$

The objective of the RRA problem is to find the optimal policy $\pi^{\text{CM}*}$ that maximizes the expected return J^{CM} , i.e.,

$$\pi^{\text{CM}*} = \arg \max_{\pi^{\text{CM}}} J^{\text{CM}}. \quad (32)$$

IV. DRL SOLUTION FOR RRA

The proposed DRL algorithm to solve the RRA problem is based on Double Deep Q Networks (DDQN) [25], which is an improved version of the classical DQN algorithm, and has been widely adopted to deal with RL tasks with discrete action space. DQN combines Q-learning with a deep neural network (DNN) $Q_i^{\text{CM}}(S_{i,(k,t)}^{\text{CM}}, a_{i,(k,t)}^{\text{CM}} \mid \theta_i)$, which acts as a function approximator and maps states to actions' value estimation [26]. Since DQN's max operator selects and evaluates action $a_{i,(k,t)}^{\text{CM}}$ using the same value, the overestimated value is more likely to be chosen, leading to over-optimistic estimates. In order to avoid this, DDQN is proposed for decoupling the selection of actions and their evaluation in the target maximization [25]. Specifically, the action $a_{i,(k,t)}^{\text{CM}}$ is chosen by the current Q networks $Q_i^{\text{CM}}(S_{i,(k,t)}^{\text{CM}}, a_{i,(k,t)}^{\text{CM}} \mid \theta_i)$ with the maximum Q value, and the evaluation of the action $a_{i,(k,t)}^{\text{CM}}$ is used by the target Q network $Q_i^{\text{CM}'}(S_{i,(k,t)}^{\text{CM}}, a_{i,(k,t)}^{\text{CM}} \mid \theta_i')$.

When tackling the RRA problem in this paper, an effective policy cannot be learned by direct application of DDQN due to the following two challenges.

Multi-agent problem Similar to PC, RRA is also a multi-agent problem. However, the simple IL approach is not suitable for RRA since the agents cooperatively optimize the *global* expected cumulative reward, and it is hard for each agent to deduce its own contribution to the global reward. The challenge in distinguishing the agents' credit makes the learning of an effective policy non-trivial.

Long-range planning with sparse reward problem

According to (22), the RRA problem involves optimization over a time horizon of KT communication intervals, where T consecutive communication intervals correspond to a control interval. In the reward function defined in (24), $R_{I,(k,t)}$ is provided to the agent in every communication interval (k, t) , while $R_{V,(k,T)}$ is only available at the last communication interval $T - 1$ of each control interval $k \in \mathcal{K}$. The sparse-reward tasks coupled with long-horizon planning become prohibitively difficult as highly specific action sequences must be executed prior to observing any nontrivial feedback.

In the following, we propose the MTCC-RRA algorithm based on DDQN, where we apply the *reward shaping* and *RBPER* [15] techniques to deal with these two challenges, respectively.

1) *Reward shaping*: To overcome the multi-agent challenge as described above, we design an individual RRA reward $R_{i,(k,t)}^{\text{CM}}$ for each vehicle $i \in \mathcal{V} \setminus \{N - 1\}$ as

$$R_{i,(k,t)}^{\text{CM}} = \begin{cases} \kappa_1 \hat{R}_{I,i,(k,t)}, & t < T - 1 \\ \kappa_1 \hat{R}_{I,i,(k,t)} + \kappa_2 A_{\hat{\pi}_{i+1}^{\text{CL}^*}} \left(x_{i+1,k+1}, \pi_{i+1}^{\text{CL}}(S_{i+1,k+1}^{\text{CL}}) \right), & t = T - 1 \end{cases}, \quad (33)$$

where

$$\hat{R}_{I,i,(k,t)} = r_{m,(k,t)} - r_{m,(k,t)}^{\bar{i}} \quad (34)$$

is the difference reward that reflects the data rate loss of V2I link m resulting from sharing the sub-channel with V2V link i , while the sub-channel allocations of the other V2V links remain the same. $r_{m,(k,t)}$ is the actual data rate of V2I link m at communication interval (k, t) , while $r_{m,(k,t)}^{\bar{i}}$ is the data rate of V2I link m if V2V link i does not occupy sub-channel m . All the channel gains for the calculation of $r_{m,(k,t)}^{\bar{i}}$ are included in the state space. Since $\hat{R}_{I,i,(k,t)}$ is a negative value, maximizing its value over time is equivalent to minimizing the throughput loss of V2I links caused by interference from V2V link i .

Moreover, since the RRA decisions at vehicle i will only affect the observation delay and thus the PC performance of its follower $i + 1$, the advantage function $A_{\hat{\pi}_{i+1}^{\text{CL}^*}} \left(x_{i+1,k+1}, \pi_{i+1}^{\text{CL}}(S_{i+1,k+1}^{\text{CL}}) \right)$ is a part of the individual reward of vehicle i at each communication interval $(k, T - 1)$ in (33).

Remark 2 (Obtaining $\hat{R}_{I,i,(k,t)}$ and $A_{\hat{\pi}_{i+1}^{\text{CL}^}}$ in (33))*: We consider the centralized training decentralized execution paradigm [27] for MTCC-RRA, so that the global state and global action are available to the RRA agents during training. Therefore, both $r_{m,(k,t)}$ and $r_{m,(k,t)}^{\bar{i}}$ can be calculated by (6) and (9), where $\theta_{i,m,(k,t)}$ in (6) for V2V link i is set to 1 and 0, respectively. Then, $\hat{R}_{I,i,(k,t)}$ can be derived by (34). Moreover, the PC advantage function $A_{\hat{\pi}_{i+1}^{\text{CL}^*}} \left(x_{i+1,k+1}, \pi_{i+1}^{\text{CL}}(S_{i+1,k+1}^{\text{CL}}) \right)$ of follower $i + 1$ is available to its predecessor i due to the centralized training assumption. Note that each RRA agent trains a local critic instead of all the agents collectively training

a centralized critic, since the global reward is decomposed into the individual rewards of the RRA agents.

2) *Reward backpropagation prioritized experience replay*: To resolve the difficulty of training for long-range planning with sparse reward, RBPER assigns a higher sampling priority β to the transition $\left(S_{i,(k,T-1)}^{\text{CM}}, a_{i,(k,T-1)}^{\text{CM}}, R_{i,(k,T-1)}^{\text{CM}}, S_{i,(k,T)}^{\text{CM}} \right)$ of the last communication interval $T - 1$ of each control interval k and then propagates the priority β back to the transition of the previous communication interval $T - 2$ once it has been sampled for training and so on. The process will keep going until β is propagated to the transition $\left(S_{i,(k,0)}^{\text{CM}}, a_{i,(k,0)}^{\text{CM}}, R_{i,(k,0)}^{\text{CM}}, S_{i,(k,1)}^{\text{CM}} \right)$ of the first communication interval 0 of each control interval k . Then priority β will be propagated back to the transition $\left(S_{i,(k,T-1)}^{\text{CM}}, a_{i,(k,T-1)}^{\text{CM}}, R_{i,(k,T-1)}^{\text{CM}}, S_{i,(k,T)}^{\text{CM}} \right)$ and start a new round of propagation. The high priority will be decayed by factor $\zeta \in (0, 1)$, once we start a new round until it goes down to normal priority 1. Thus, the target Q, $Q_i^{\text{CM}'}(S_{i,(k,t+1)}^{\text{CM}}, a_{i,(k,t+1)}^{\text{CM}} | \theta_i')$, of the previous communication interval t will be established on the basis of the trained Q, $Q_i^{\text{CM}}(S_{i,(k,t+1)}^{\text{CM}}, a_{i,(k,t+1)}^{\text{CM}} | \theta_i)$, of the next communication interval $t + 1$, improving the accuracy of Q value estimation and training efficiency.

V. JOINT OPTIMIZATION OF MULTI-TIMESCALE CONTROL AND COMMUNICATIONS

As mentioned above, it is assumed that the RRA (resp. PC) policy is available when training the PC (resp. RRA) policy. Due to the interactions between RRA and PC, it is not possible to learn the optimal RRA or PC policy first without knowing the other optimal policy. To solve this dilemma, the most straightforward approach is to train the proposed MTCC-PC and MTCC-RRA algorithms simultaneously. However, training both algorithms from scratch poses a great challenge in convergence, especially for RRA whose rewards depend on the PC. The rewards calculated from a PC model that has not reached convergence yet are misleading to RRA, resulting in substantial difficulty in learning an efficient RRA policy. Moreover, it is more challenging for MTCC-RRA algorithm to reach convergence than MTCC-PC algorithm under the multi-timescale framework, since the number of communication intervals is T times that of the control intervals per episode. In this paper, we propose to iteratively train MTCC-PC and MTCC-RRA algorithms, where each iteration $z \in \{1, \dots, Z\}$ consists of the following two steps.

Step 1 The MTCC-PC is trained for E^{CL} episodes under the observation delay induced by the RRA policy learned in Step 2 of iteration $z - 1$. The RRA transitions $\left(S_{i,(k,t)}^{\text{CM}}, a_{i,(k,t)}^{\text{CM}}, R_{i,(k,t)}^{\text{CM}}, S_{i,(k,t+1)}^{\text{CM}} \right)$ are stored in an experience replay buffer D_i^{CM} for all episodes $E > E^{\text{CLTh}}$.

Step 2 The MTCC-RRA is trained for E^{CM} episodes, where the PC model learned in Step 1 of iteration z is used to calculate vehicle trajectories and RRA rewards. The PC transitions $\left(S_{i,k}^{\text{CL}}, a_{i,k}^{\text{CL}}, R_{i,k}^{\text{CL}}, S_{i,k}^{\text{CL}} \right)$ are stored in an experience replay buffer D_i^{CL} for all episodes $E > E^{\text{CMTh}}$.

Note that it is essential to train the MTCC-PC first in Step 1, and the MTCC-RRA next in Step 2. This is because RRA training heavily depends on the PC advantage function to calculate reasonable reward signals as given in (28). The initial RRA policy randomly selects the RRA actions according to uniform distribution.

To improve the sample efficiency, the RRA (resp. PC) experience is stored in the experience replay buffer D_i^{CM} (resp. D_i^{CL}) when training the PC (resp. RRA) model. The stored experience can be leveraged at the initial phase of training when RRA (resp. PC) is trained in Step 2 (resp. Step 1 of the next iteration). To ensure that the experience are generated by relatively well-learned PC (resp. RRA) policies in the previous step, only the transitions at the later stage of PC (resp. RRA) training when $E > E^{\text{CL}_{\text{Th}}}$ (resp. $E > E^{\text{CM}_{\text{Th}}}$) are used.

Although it takes multiple iterations to obtain optimal results in MTCC, we will demonstrate in Section VI that satisfactory policies can be learned in only one iteration.

The proposed MTCC training algorithm is given in Algorithm 1, where the MTCC-PC and MTCC-RRA algorithms are iteratively trained for Z times. In each iteration, the function $\text{MTCC}(\text{PC_TRAIN} = 1, \text{RRA_TRAIN} = 0, E^{\text{CL}}, \text{ITER} = z)$ corresponds to the MTCC-PC algorithm and the function $\text{MTCC}(\text{PC_TRAIN} = 0, \text{RRA_TRAIN} = 1, E^{\text{CM}}, \text{ITER} = z)$ corresponds to the MTCC-RRA algorithm. To improve the sample efficiency, the transitions at the later stage of PC (resp. RRA) training when $E > E^{\text{CL}_{\text{Th}}}$ (resp. $E > E^{\text{CM}_{\text{Th}}}$) are used at the initial phase of training when RRA (resp. PC) is trained in Step 2 (resp. Step 1 of the next iteration). The pseudocode of the function MTCC is provided in Funtion 1. When the input argument $\text{PC_TRAIN} = 1$, the DDPG algorithm is used to train the PC agent. When the input argument $\text{RRA_TRAIN} = 1$, the DDQN-RBPER algorithm is used to train the RRA agent. The pseudocode of the DDPG algorithm and DQN-RBPER algorithm can be found in [16] and [15], respectively. The DDQN-RBPER algorithm extends DQN-RBPER with DDQN [25].

The unified DRL training framework for MTCC is illustrated in Fig. 1. Each vehicle i contains a PC agent and an RRA agent, except for the leading vehicle 0 which only has an RRA agent and the last vehicle $N - 1$ which only has a PC agent. The environment can also be divided into the PC environment, which provides driving status and PC rewards to the PC agents; and the RRA environment, which provides wireless channel and queue states as well as V2I-related rewards to the RRA agents. Meanwhile, the PC (resp. RRA) agents also interact with the RRA (resp. PC) environment/agents. For each RRA agent of vehicle $i \in \mathcal{V} \setminus \{N - 1\}$, the PC advantage function $A_{\pi_{i+1,k+1}^{\text{CL}*}}$ of vehicle $i + 1$ at control interval $k + 1$ is used to calculate the PC-related portion of RRA reward $R_{i,(k,T-1)}^{\text{CM}}$ at communication interval $(k, T - 1)$. Moreover, the RRA state $S_{i,(k,t)}^{\text{CM}}$ of vehicle i at communication interval (k, t) is augmented with its PC action history $\{a_{i,k'}^{\text{CL}}\}_{k'=k-\tau_{\text{max}}+1}^k$. Meanwhile, for each PC agent of vehicle $i + 1 \in \mathcal{V} \setminus \{0\}$, its state $S_{i+1,k+1}^{\text{CL}}$ at control interval $k + 1$ depends on the observation delay $\tau_{i+1,k+1}$, which is provided by the RRA environment and determined by the queue state $q_{i,(k,T)}^{\text{CAM}}$ of

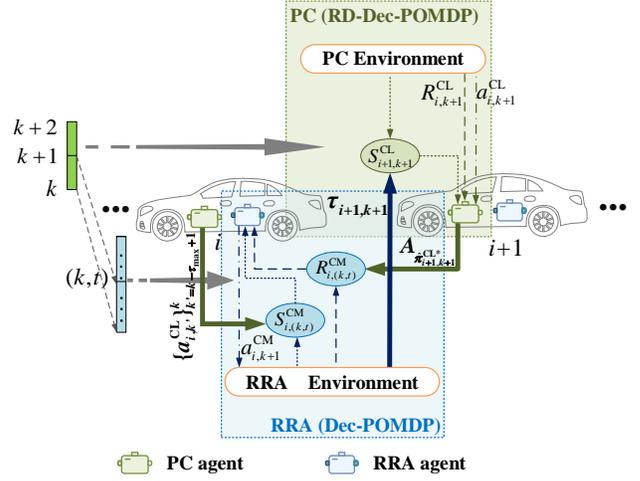


Fig. 1. The framework of MTCC algorithm

vehicle i at communication interval (k, T) .

After the PC agents and RRA agents are trained in the centralized fashion, each vehicle makes PC and/or RRA decisions in a decentralized manner based on its local observations of the PC and RRA states. The decisions can be made in real-time since only the forward propagation in the DNNs is involved during execution.

Algorithm 1 MTCC algorithm

- 1: **for** follower $i \in \mathcal{V} \setminus \{0\}$ **do**
 - 2: Initialize PC actor network $\mu_i^{\text{CL}}(s|\theta_i^\mu)$ and critic network $Q_i^{\text{CL}}(s, u|\theta_i^Q)$ with random weights θ_i^μ and θ_i^Q
 - 3: Initialize PC target network weights with $\theta_i^{\mu'} = \theta_i^\mu$ and $\theta_i^{Q'} = \theta_i^Q$
 - 4: Initialize replay buffer D_i^{CL}
 - 5: **end for**
 - 6: **for** predecessor $i \in \mathcal{V} \setminus \{N - 1\}$ **do**
 - 7: Initialize RRA Q networks $Q_i^{\text{CM}}(s, a|\theta_i)$ with random weights θ_i
 - 8: Initialize RRA target network weights with $\theta_i' = \theta_i$
 - 9: Initialize replay buffer D_i^{CM}
 - 10: **end for**
 - 11: **for** Iteration $z = 1, \dots, Z$ **do**
 - 12: **Step 1:**
 - 13: Function $\text{MTCC}(\text{PC_TRAIN} = 1, \text{RRA_TRAIN} = 0, E^{\text{CL}}, \text{ITER} = z)$
 - 14: Delete RRA experience for episodes $E \leq E^{\text{CL}_{\text{Th}}}$ from D_i^{CM}
 - 15: **Step 2:**
 - 16: Function $\text{MTCC}(\text{PC_TRAIN} = 0, \text{RRA_TRAIN} = 1, E^{\text{CM}}, \text{ITER} = z)$
 - 17: Delete PC experience for episodes $E \leq E^{\text{CM}_{\text{Th}}}$ from D_i^{CL}
 - 18: **end for**
-

Function 1 MTCC(PC_TRAIN, RRA_TRAIN, E , ITER)

```

1: for episode  $e = 1, \dots, E$  do
2:   for control interval  $k = 0, \dots, K - 1$  do
3:     for communication interval  $t = 0, \dots, T - 1$  do
4:       for predecessor  $i \in \mathcal{V} \setminus \{N - 1\}$  do
5:         Update the vehicle position and receive RRA
           state  $S_{i,(k,t)}^{\text{CM}}$ 
6:         if RRA_TRAIN = 1 then
7:           Select RRA action  $a_{i,(k,t)}^{\text{CM}}$  according
           to the  $\epsilon$ -greedy policy with respect to
            $Q_i^{\text{CM}}(S_{i,(k,t)}^{\text{CM}}, a|\theta_i)$ 
8:         else
9:           if ITER = 1 then
10:            Select a random RRA action  $a_{i,(k,t)}^{\text{CM}}$ 
11:          else
12:            Select RRA action  $a_{i,(k,t)}^{\text{CM}}$  according
            to the greedy policy with respect to
             $Q_i^{\text{CM}}(S_{i,(k,t)}^{\text{CM}}, a|\theta_i)$ 
13:          end if
14:        end if
15:        Execute RRA action  $a_{i,(k,t)}^{\text{CM}}$  and observe RRA
           reward  $R_{i,(k,t)}^{\text{CM}}$  and next RRA state  $S_{i,(k,t+1)}^{\text{CM}}$ 
16:        Store RRA transition  $(S_{i,(k,t)}^{\text{CM}}, a_{i,(k,t)}^{\text{CM}}, R_{i,(k,t)}^{\text{CM}},$ 
            $S_{i,(k,t+1)}^{\text{CM}})$  with priority  $p$  in  $D_i^{\text{CM}}$ 
17:        if RRA_TRAIN = 1 then
18:          Update  $\theta_i, \theta'_i$  with DDQN-RBPER algorithm
19:        end if
20:      end for
21:    end for
22:    for follower  $i \in \mathcal{V} \setminus \{0\}$  do
23:      Calculate observation delay  $\tau_{i,k}$  and receive PC
           state  $S_{i,k}^{\text{CL}}$ 
24:      Select PC action  $a_{i,k}^{\text{CL}} = \mu_i^{\text{CL}}(S_{i,k}^{\text{CL}}|\theta_i^\mu)$ 
25:      if PC_TRAIN = 1 then
26:        Add exploration noise to PC action  $a_{i,k}^{\text{CL}} =$ 
            $a_{i,k}^{\text{CL}} + \mathcal{N}_t$ 
27:      end if
28:      Store PC transition  $(S_{i,k}^{\text{CL}}, a_{i,k}^{\text{CL}}, R_{i,k}^{\text{CL}}, S_{i,k+1}^{\text{CL}})$  in
            $D_i^{\text{CL}}$ 
29:      if PC_TRAIN = 1 then
30:        Update  $\theta_i^\mu, \theta_i^Q, \theta_i^{\mu'}, \theta_i^{Q'}$  with DDPG algorithm
31:      end if
32:    end for
33:  end for
34: end for

```

VI. EXPERIMENTAL RESULTS

In this section, we design experiments to answer the following questions: (1) Can MTCC-RRA outperform the state-of-the-art control-aware communications by capturing ‘‘VoI per control interval’’ through its state space and reward design? (2) Can the reward shaping and RBPER mechanisms efficiently tackle the multi-agent and sparse reward problems for MTCC-RRA? (3) Can the satisfactory policies for MTCC be learned in only one iteration? The experimental results of the proposed MTCC algorithm and the baseline DRL algorithms are

TABLE I
HYPER-PARAMETERS OF THE DRL ALGORITHMS FOR TRAINING

Parameter	Value
Network size	256, 64
Activation function	relu, relu, relu, linear
Learning rate	0.0001
Exploration rate ϵ	$1 \rightarrow 0.05$
Higher sampling priority β	100
Decay factor ζ	0.2
Batch size N_b	64
Replay buffer size	200000
Reward discount factor γ	0.99979
Target update frequency of DDQN N^-	4
Final layer weights/biases initialization	Random uniform distribution $[-3 \times 10^{-3}, 3 \times 10^{-3}]$
Other layer weights/biases initialization	Random uniform distribution $[-\frac{1}{\sqrt{f}}, \frac{1}{\sqrt{f}}]$ (f is the fan-in of the layer)

presented and compared for the above purposes.

A. Experimental Setup

We conduct experiments based on the PC and RRA environments provided in Part I of this two-part paper. In addition to the hyper-parameters of DDPG for solving the communication-aware PC problem, we summarized the hyper-parameters for DDQN solving the control-aware RRA problem in Table I. The values of all the hyper-parameters were selected by performing a grid search as in [26], using the values reported in [25] as a reference. DDQN has two hidden layers with 256 and 64 nodes, respectively, where 128 of the nodes in the first hidden layer are modeled as Long Short-Term Memory (LSTM) for processing PC action history $\{a_{i,k'}^{\text{CL}}\}_{k'=k-\tau_{\max}+1}^k$ in $S_{i,(k,t)}^{\text{CM}}$ and the rest of the 128 nodes in the first hidden layer are used for the remaining dimensions of state $S_{i,(k,t)}^{\text{CM}}$. Then all the 256 nodes are connected to the second layer with 64 nodes. The target network of DDQN is updated every 4 communication intervals.

1) Baseline algorithms:

- **Delay-aware RRA (Delay-RRA)**, which is an improved version of the DRL-based RRA algorithm for C-V2X communications in [9], [11], [12]. [9] focuses on RRA within one control interval and aims to minimize the latency and maximize the success probability of delivering the CAM generated at the beginning of the control interval. To deal with the sparse reward problem, the V2V transmission rate at each time step is included in the reward function to expedite convergence. Since [9] does not address the credit assignment problem in MARL, we improve the algorithm by applying reward shaping similar to MTCC for a fair comparison. Moreover, to apply the algorithm in a longer time horizon with multiple control intervals, we assume that an agent replaces any old CAM that has not yet been fully delivered in the previous control interval with the newly generated CAM

at the beginning of each control interval. Therefore, the reward function $R_{i,(k,t)}^{\text{IR}}$ of Delay-RRA is given as

$$R_{i,(k,t)}^{\text{IR}} = \lambda_1 \hat{R}_{1,i,(k,t)} + \lambda_2 (r_{i,(k,t)}|_{q_{i,(k,t)}^{\text{CAM}} > 0} + G|_{q_{i,(k,t)}^{\text{CAM}} = 0}), \quad (35)$$

where $\lambda_1 = 0.001/W$ and $\lambda_2 = 0.1/W$. Notice that G is a tuned hyper-parameter encouraging early completion of transmission for CAM whose value is greater than the largest $r_{i,(k,t)}$ ever obtained. We set $G = 10W$ in the experiments.

- **AoI-aware RRA (AoI-RRA)**, which is similar to the DRL-based RRA in MTCC except that it aims to minimize cumulative AoI instead of maximizing cumulative ‘‘VoI per control interval’’. Specifically, the advantage function $A_{\hat{\pi}_{i+1}^{\text{CL}^*}}(x_{i+1,k+1}, \pi_{i+1}^{\text{CL}}(S_{i+1,k+1}^{\text{CL}}))$ in (33) is replaced by $-AoI_{i+1,k}$, where $AoI_{i+1,k}$ is the AoI of vehicle $i + 1$ that evolves according to

$$AoI_{i+1,k+1} = \begin{cases} T & \text{if } q_{i,(k,T)}^{\text{CAM}} = 0 \\ AoI_{i+1,k} + T & \text{if } q_{i,(k,T)}^{\text{CAM}} > 0 \end{cases}. \quad (36)$$

For Delay-RRA and AoI-RRA, $q_{i,(k,0)}^{\text{CAM}}$ is always 1 and no longer equals to $q_{i,(k-1,T)}^{\text{CAM}}$.

- **MTCC without reward shaping (MTCC_wo_RS)**, where the global reward in (24) is used in MTCC for training DRL-based RRA, instead of the individual reward in (33) after reward shaping.
- **MTCC without RBPER (MTCC_wo_RBPER)**, where DDQN without RBPER is used in MTCC for training DRL-based RRA.

2) *Experiment design*: To answer the three questions mentioned at the beginning of Section VI, we design the following three experiments. The proposed MTCC is trained in one iteration for Experiments (1) and (2).

- **Experiment (1): Performance comparison of MTCC-RRA with Delay-RRA and AoI-RRA.** To answer Question (1), all three algorithms are trained for $E^{\text{CM}} = 100$ episodes. The trained PC policy by MTCC-PC is adopted during both the training and testing of the RRA algorithms. Therefore, the PC performance difference between Delay-RRA, AoI-RRA, and MTCC-RRA is only due to the difference in induced observation delay by the RRA algorithms.
- **Experiment (2): Performance comparison of MTCC-RRA with MTCC_wo_RS and MTCC_wo_RBPER.** To answer Question (2), MTCC_wo_RS and MTCC_wo_RBPER are trained for $E^{\text{CM}} = 100$ episodes and the resulting performance is compared with that of MTCC-RRA.
- **Experiment (3): Performance Comparison of MTCC in one and two iterations.** To answer Question (3), we train MTCC for a second iteration and observe the resultant performance gain.

Note that the episode thresholds for storing experience are set to $E^{\text{CL}_{\text{Th}}} = \frac{E^{\text{CL}}}{5}$ and $E^{\text{CM}_{\text{Th}}} = \frac{E^{\text{CM}}}{5}$, respectively.

B. Performance Comparison of MTCC-RRA, Delay-RRA, and AoI-RRA.

1) *Performance for testing data*: The RRA performance are reported in Table II for MTCC-RRA, Delay-RRA, and AoI-RRA, respectively, which are obtained by averaging the RRA return over 100 testing episodes when training is completed. The RRA return is derived as the cumulative global RRA reward defined in (24) overall communication intervals in one control episode. Since the global RRA reward consists of two parts, i.e., the sum throughput of all V2I links and the sum PC performance degradation due to delayed observation for all followers, we also report the average sum V2I throughput and sum PC performance in Table II. The former is obtained by averaging the sum throughput for all V2I links over 100 testing episodes, while the latter is obtained by averaging the sum PC returns of all followers. Note that the sum PC performance and the sum PC performance degradation are essentially the same since their values are different by the same fixed amount for all the algorithms, where the difference corresponds to the optimal sum PC performance without observation delay.

As shown in Table II, MTCC-RRA outperforms AoI-RRA and Delay-RRA algorithms by 58.44% and 55.81%, respectively, in terms of RRA performance. The result is not surprising since MTCC-RRA aims at maximizing the RRA performance corresponding to the cumulative global RRA reward in (24), while Delay-RRA and AoI-RRA have different reward functions.

The reward function of MTCC-RRA differs from those of Delay-RRA and AoI-RRA in that MTCC-RRA aims at minimizing the PC performance degradation due to observation delay instead of minimizing the delay or AoI, while all the reward functions also target at maximizing the V2I throughput. Therefore, MTCC-RRA can achieve larger V2I throughput by allocating more radio resources to V2I links when larger observation delays do not have significant impacts on the PC performance. This is further proved by a close examination of the average sum V2I throughput and sum PC performance in Table II. In terms of the average sum V2I throughput, MTCC-RRA outperforms AoI-RRA and Delay-RRA algorithms by 20.02% and 20.72%, respectively. Meanwhile, MTCC-RRA is slightly inferior to AoI-RRA and Delay-RRA in terms of sum PC performance by 8.14%. Therefore, MTCC-RRA can strike a better trade-off between PC performance and V2I throughput compared with the other two RRA algorithms.

2) *Testing results of one episode*: The key reasons for the superior performance of MTCC-RRA are quantifying the ‘‘VoI per control interval’’ in the reward function and including the control action history in the RRA state. To better understand the positive effects of these design factors, we zoom in to one control episode, where the tracking errors $e_{pi,k}$ and $e_{vi,k}$ of each follower $i \in \mathcal{V} \setminus \{0\}$ as well as the acceleration $acc_{i,k}$ and control input $a_{i,k}^{\text{CL}}$ of each vehicle $i \in \mathcal{V}$ for all time steps $k \in \{1, 2, \dots, 120\}$ are plotted in Fig. 2 for MTCC-RRA, Delay-RRA, and AoI-RRA. In addition, the observation delay, $\tau_{i,k}$, of each follower $i \in \mathcal{V} \setminus \{0\}$ for each time step k are also plotted.

It can be observed that the performance curves of Delay-

TABLE II
PERFORMANCE AFTER TRAINING RRA IN ONE ITERATION. WE PRESENT RRA PERFORMANCE AND AVERAGE SUM V2I THROUGHPUT, AND SUM PC PERFORMANCE FOR MTCC-RRA, DELAY-RRA, AND AOI-RRA, MTCC_wo_RS, AND MTCC_wo_RBPER, RESPECTIVELY.

	MTCC-RRA	Delay-RRA	AoI-RRA	MTCC_wo_RS	MTCC_wo_RBPER
RRA performance	68.2553	43.0803	43.8063	58.3390	-209.1657
Average sum V2I throughput (Mbps)	4.3982	3.6646	3.6432	4.0040	4.7523
Sum PC performance	-1.5390	-1.4231	-1.4231	-3.2671	-28.1860

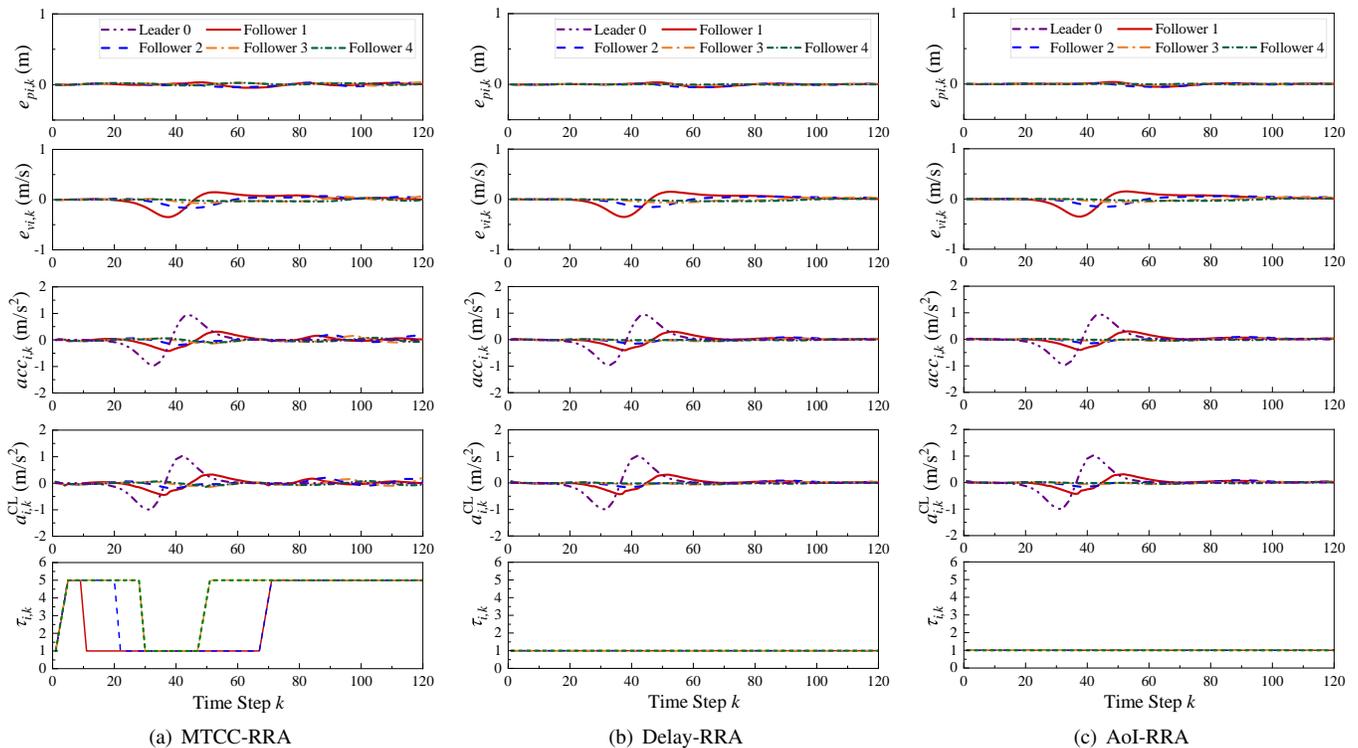


Fig. 2. Results of a specific test control episode. The tracking errors $e_{pi,k}$, $e_{vi,k}$, and $acc_{i,k}$ along with the control input $a_{i,k}^{CL}$ and observation delay $\tau_{i,k}$ of each follower i are represented as different curves, respectively.

RRA and AoI-RRA are the same since the minimum observation delay of 1 control interval is always achieved for both algorithms. Meanwhile, the overall shapes of $e_{pi,k}$, $e_{vi,k}$, $acc_{i,k}$, and $a_{i,k}^{CL}$ curves for MTCC-RRA look very similar to those of Delay-RRA and AoI-RRA, while the $\tau_{i,k}$ curves are quite different. Instead of achieving a constant minimum delay at all time steps, it can be seen that the observation delay of MTCC-RRA for each follower $i \in \mathcal{V} \setminus \{0\}$ is minimum only when the control input $a_{i-1,k}^{CL}$ of its predecessor $i-1$ changes significantly and rapidly. When $a_{i-1,k}^{CL}$ remains relatively constant, the follower i generally has large observation delays, which means that the predecessor $i-1$ refrains from sharing the V2I sub-channels for CAM transmission to support larger V2I throughput. This observation demonstrates that the MTCC-RRA agent has the ability to allocate radio resources based on the VoI per control interval.

As a result, the PC performance of MTCC-RRA is only slightly affected compared to those of Delay-RRA and AoI-RRA. Fig. 2 shows that the convergence speed to steady state for all algorithms are similar, and the PC performance differences of the algorithms are reflected in the oscillations of the tracking errors, acceleration, and control input. After

the steady state is reached and the tracking errors are almost zero, the performance curves of MTCC-RRA have slightly larger oscillations compared with those in Delay-RRA and AoI-RRA. As the oscillations are negligible, MTCC-RRA can still ensure driving safety and comfort. Moreover, the amplitudes of the oscillations in $e_{pi,k}$, $e_{vi,k}$, and $acc_{i,k}$ for each follower $i \in \{1, 2\}$ are smaller than those of their respective predecessors $i-1$ for all algorithms, which means that the string stability of the platoons can be achieved by MTCC-RRA despite the larger observation delays.

3) *Convergence properties*: The RRA performance of Delay-RRA, AoI-RRA, and MTCC-RRA are evaluated periodically during training by testing in an environment where the trained MTCC-PC policy is used to calculate vehicle trajectories and RRA rewards. Specifically, we run 1 test control episodes after every 1 training control episode and calculate the RRA return as the performance for the latest 1 training episode. The RRA performance as a function of the number of training episodes with MTCC-RRA, Delay-RRA, and AoI-RRA algorithms is plotted in Fig. 3. It is shown that the convergence rate for all three algorithms are similar, with a ranking of MTCC-RRA > AoI-RRA > Delay-RRA.

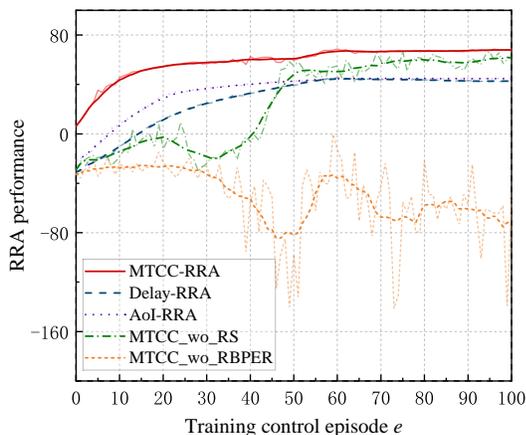


Fig. 3. RRA Performance during DRL-based RRA algorithms training in one iteration. The vertical axis corresponds to the returns over 1 test episode. The dark curves correspond to smoothed curves and the light color curves correspond to the original curves.

C. Performance Comparison of MTCC-RRA, MTCC_wo_RS, and MTCC_wo_RBPER.

The RRA performance, average sum V2I throughput, and sum PC performance are reported in Table II for MTCC_wo_RS and MTCC_wo_RBPER, respectively. In terms of RRA performance, MTCC-RRA outperforms MTCC_wo_RS and MTCC_wo_RBPER algorithms by 16.95% and 132.63%, respectively. The superior performance of MTCC-RRA over MTCC_wo_RS demonstrates the effectiveness of the reward shaping technique in efficiently distinguishing the multiple agents' credit in the global reward. Moreover, the significant performance gain of MTCC-RRA over MTCC_wo_RBPER shows that the RBPER technique is essential in solving the long-range planning with sparse reward problem. The same conclusion can also be drawn from the convergence curves of MTCC_wo_RS and MTCC_wo_RBPER shown in Fig. 3. Both performance curves have large fluctuations, and that of MTCC_wo_RBPER is not convergent.

D. Performance Comparison of MTCC in one and two iterations

In the above experiments, MTCC is only trained in one iteration, i.e., $Z = 1$, in Algorithm 1. In the following, we train MTCC for one additional iteration to observe the resultant performance gain. For the second iteration, the performance curves for the PC policy, i.e., MTCC-PC, during training in Step 1 and for the RRA policy, i.e., MTCC-RRA, during training in Step 2 are given in Fig. 4, respectively. The performance curve of MTCC-PC is obtained in a similar way as that in Fig.3 in Part I, except that the MTCC-PC policy is trained and tested under a delayed environment generated by the MTCC-RRA policy learned in Step 2 of iteration 1. Similarly, the performance curve of MTCC-RRA is obtained in a similar way as that in Fig. 3, except that the MTCC-RRA policy is trained and tested under the PC environment generated by the MTCC-PC policy learned in Step 1 of iteration 2. From Fig. 4, it can be observed that the sum PC performance and RRA performance improve with

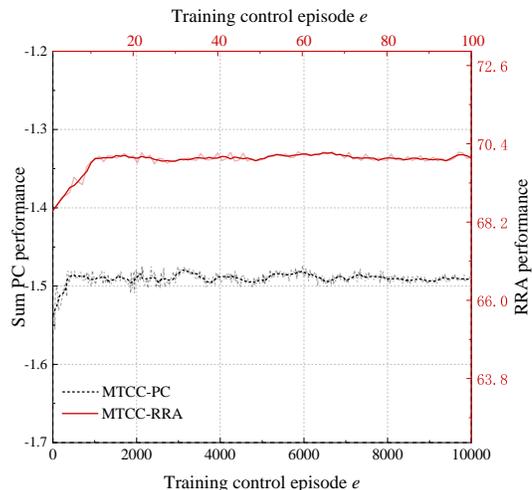


Fig. 4. Sum PC performance and RRA performance during training in the second iteration.

increasing training episodes in Step 1 and Step 2 of the second iteration, respectively. However, the performance gains are not significant.

Specifically, the sum PC performance of MTCC-PC is around -1.55 at the beginning of training and finally converges around -1.49 towards the end, with a performance gain of 3.87%. The performance gain in Step 1 is due to the fact that the MTCC-PC policy is trained under the observation delay generated by the de facto MTCC-RRA policy instead of the random RRA policy. Similarly, the RRA performance of MTCC-RRA is around 68.4 at the beginning of training and finally converges around 70.2 towards the end, with a performance gain of 2.63%. The performance gain in Step 2 is because the de facto MTCC-PC policy trained in iteration 2 is adopted instead of that trained in iteration 1 to simulate the PC environment, so that the reward function of MTCC-RRA calculated from the MTCC-PC advantage function can more accurately represent the PC performance degradation due to observation delay.

The above results show that satisfactory policies for MTCC can be learned in one iteration and a second iteration is merely a fine-tuning for small performance improvement. Moreover, it can be anticipated that the performance gain will keep decreasing when training in more iterations.

VII. CONCLUSION

In this pair of papers, we have studied how to solve the joint optimization of the MTCC problem in the C-V2X system. We have established a unified DRL framework for MTCC and proposed the MTCC-PC algorithm to solve the communication-aware PC sub-problem in Part I of this two-part paper. In this paper, we have proposed the MTCC-RRA algorithm to solve the control-aware RRA sub-problem. Specifically, the RRA decisions have been made based on the advantage function of the PC model, which provides a fine-grained VoI per control interval. The multi-agent problem

and the sparse reward problem in RRA have been tackled with the reward shaping and RBPER techniques, respectively. Finally, we have designed a sample- and computational-efficient training approach to jointly learn the PC and RRA policies iteratively. The experimental results have shown that the proposed MTCC-RRA algorithm outperforms the baseline DRL algorithms, striking a good trade-off between minimizing PC performance degradation due to observation delay and maximizing V2I throughput. Moreover, it is demonstrated that the joint training approach can learn satisfactory PC and RRA policies in only one iteration. In our future work, we will apply the proposed MTCC framework to other autonomous driving tasks such as lane changing and merging.

APPENDIX

A. Proof of Lemma 1

As explained in Section II.D, we convert the finite-horizon PC problem to an infinite-horizon one by calculating the TD target $y_{i,K-1}$ using (19) for the last control interval $K-1$ instead of setting $y_{i,K-1}$ to be the immediate reward. Therefore, in the following proof, we replace K with ∞ on the RHS of (26).

$$\begin{aligned}
& \mathbb{E}_{\pi_{\text{CM}}}\mathbb{E}_{\pi_{\text{CL}}}\left[\sum_{k=0}^{\infty}\gamma^k A_{\hat{\pi}_i^{\text{CL}^*}}(x_{i,k},\pi_i^{\text{CL}}(S_{i,k}^{\text{CL}}))\right] \\
& \stackrel{(a)}{=} \mathbb{E}_{\pi_{\text{CM}}}\mathbb{E}_{\pi_{\text{CL}}}\left[\sum_{k=0}^{\infty}\gamma^k\left(Q_{\hat{\pi}_i^{\text{CL}^*}}(x_{i,k},\pi_i^{\text{CL}}(S_{i,k}^{\text{CL}}))\right.\right. \\
& \quad \left.\left.-V_{\hat{\pi}_i^{\text{CL}^*}}(x_{i,k})\right)\right] \\
& \stackrel{(b)}{=} \mathbb{E}_{\pi_{\text{CM}}}\mathbb{E}_{\pi_{\text{CL}}}\left[\sum_{k=0}^{\infty}\gamma^k\left(\mathbb{E}_{x_{i,k+1}}[R_{i,k}^{\text{CL}}+\gamma V_{\hat{\pi}_i^{\text{CL}^*}}(x_{i,k+1})]\right.\right. \\
& \quad \left.\left.-V_{\hat{\pi}_i^{\text{CL}^*}}(x_{i,k})\right)\right] \\
& \stackrel{(c)}{=} \mathbb{E}_{\pi_{\text{CM}}}\mathbb{E}_{\pi_{\text{CL}}}\left[\sum_{k=0}^{\infty}\gamma^k\left(R_{i,k}^{\text{CL}}+\gamma V_{\hat{\pi}_i^{\text{CL}^*}}(x_{i,k+1})\right.\right. \\
& \quad \left.\left.-V_{\hat{\pi}_i^{\text{CL}^*}}(x_{i,k})\right)\right] \\
& = \mathbb{E}_{\pi_{\text{CM}}}\mathbb{E}_{\pi_{\text{CL}}}\left[\sum_{k=0}^{\infty}\gamma^k R_{i,k}^{\text{CL}}-V_{\hat{\pi}_i^{\text{CL}^*}}(x_{i,0})\right] \\
& = \mathbb{E}_{\pi_{\text{CM}}}\mathbb{E}_{\pi_{\text{CL}}}\left[\sum_{k=0}^{\infty}\gamma^k R_{i,k}^{\text{CL}}\right]-\mathbb{E}_{x_{i,0}}\left[V_{\hat{\pi}_i^{\text{CL}^*}}(x_{i,0})\right] \\
& \stackrel{(d)}{=} J_i^{\text{CL}}-\hat{J}_i^{\text{CL}^*} \tag{37}
\end{aligned}$$

where (a) follows from (27); (b) follows from the Bellman equation; (c) follows by merging the expectation $\mathbb{E}_{x_{i,k+1}}$ into the expectation $\mathbb{E}_{\pi_{\text{CM}}}\mathbb{E}_{\pi_{\text{CL}}}$; (d) follows from (16), (21) and the definition of the value function.

B. Proof of Theorem 1

$$J^{\text{CM}} = \mathbb{E}_{\pi_{\text{CM}}}\mathbb{E}_{\pi_{\text{CL}}}\left[\sum_{k=0}^{K-1}\sum_{t=0}^{T-1}\eta^{kT+t}R_{(k,t)}^{\text{CM}}\right]$$

$$\begin{aligned}
& \stackrel{(a)}{=} \kappa_1 \sum_{m=0}^{M-1} \mathbb{E}_{\pi_{\text{CM}}}\left[\sum_{k=0}^{K-1}\sum_{t=0}^{T-1}\eta^{kT+t}r_{m,(k,t)}\right] + \kappa_2 \sum_{i=0}^{N-2} \left(\mathbb{E}_{\pi_{\text{CM}}}\right. \\
& \quad \left.\mathbb{E}_{\pi_{i+1}^{\text{CL}}}\left[\sum_{k=0}^{K-1}\eta^{kT+T-1}A_{\hat{\pi}_{i+1}^{\text{CL}^*}}(x_{i+1,k+1},\pi_{i+1}^{\text{CL}}(S_{i+1,k+1}^{\text{CL}}))\right]\right) \\
& \stackrel{(b)}{=} \kappa_1 \sum_{m=0}^{M-1} \mathbb{E}_{\pi_{\text{CM}}}\left[\sum_{k=0}^{K-1}\sum_{t=0}^{T-1}\eta^{kT+t}r_{m,(k,t)}\right] + \kappa_2 \sum_{i=0}^{N-2} \left(\mathbb{E}_{\pi_{\text{CM}}}\right. \\
& \quad \left.\mathbb{E}_{\pi_{i+1}^{\text{CL}}}\left[\sum_{k=0}^{K-1}\gamma^k A_{\hat{\pi}_{i+1}^{\text{CL}^*}}(x_{i+1,k+1},\pi_{i+1}^{\text{CL}}(S_{i+1,k+1}^{\text{CL}}))\right]\right) \\
& \stackrel{(c)}{=} \kappa_1 \sum_{m=0}^{M-1} \mathbb{E}_{\pi_{\text{CM}}}\left[\sum_{k=0}^{K-1}\sum_{t=0}^{T-1}\eta^{kT+t}r_{m,(k,t)}\right] + \kappa_2 \sum_{i=0}^{N-2} \left(J_{i+1}^{\text{CL}}-\right. \\
& \quad \left.\hat{J}_{i+1}^{\text{CL}^*}\right) \\
& \stackrel{(d)}{=} \kappa_1 \sum_{m=0}^{M-1} \mathbb{E}_{\pi_{\text{CM}}}\left[\sum_{k=0}^{K-1}\sum_{t=0}^{T-1}\eta^{kT+t}r_{m,(k,t)}\right] + \kappa_2 \Delta J^{\text{CL}}. \tag{38}
\end{aligned}$$

where (a) follows from (24) and (25); (b) holds if $\eta^{kT+T-1} = \gamma^k$; (c) follows from Lemma 1; and (d) follows from (20). Thus it is proved that the expected return as calculated in (23) is the same as the optimization objective in (22).

In the above derivation, we assume that $\eta^{kT+T-1} = \gamma^k$. Since the PC problem is essentially an infinite horizon problem as discussed in Section II.D, we have

$$\eta = \gamma^{\frac{k}{(k+1)T-1}} = \lim_{k \rightarrow +\infty} \gamma^{\frac{k}{(k+1)T-1}} = \gamma^{\frac{1}{T}}. \tag{39}$$

C. Proof of Lemma 2

If $x_{i,k}$ can be fully determined by $S_{i,k}^{\text{CL}}$ with $x_{i,k} = f(S_{i,k}^{\text{CL}})$, the relationship between $\hat{\pi}_i^{\text{CL}^*}(x_{i,k})$ and $\pi_i^{\text{CL}^*}(S_{i,k}^{\text{CL}})$ can be established as

$$a_{i,k}^{\text{CL}^*} = \hat{\pi}_i^{\text{CL}^*}(x_{i,k}) = \hat{\pi}_i^{\text{CL}^*}(f(S_{i,k}^{\text{CL}})) = \pi_i^{\text{CL}^*}(S_{i,k}^{\text{CL}}) \tag{40}$$

D. Proof of Lemma 3

The system dynamics evolve in discrete time on the basis of forward Euler discretization is

$$x_{i,k+1} = g(x_{i,k}, a_{i,k}^{\text{CL}}, a_{i-1,k}^{\text{CL}}) = Ax_{i,k} + Ba_{i,k}^{\text{CL}} + Ca_{i-1,k}^{\text{CL}}, \tag{41}$$

where

$$A = \begin{bmatrix} 1 & T & -h_i T & 0 \\ 0 & 1 & -T & T \\ 0 & 0 & 1 - \frac{T}{\tau} & 0 \\ 0 & 0 & 0 & 1 - \frac{T}{\tau} \end{bmatrix}, B = \begin{bmatrix} 0 \\ 0 \\ \frac{T}{\tau} \\ 0 \end{bmatrix}, C = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \frac{T}{\tau} \end{bmatrix}. \tag{42}$$

Then,

$$\begin{aligned}
x_{i,k} & = g(x_{i,k-1}, a_{i,k-1}^{\text{CL}}, a_{i-1,k-1}^{\text{CL}}) \\
& = g(g(x_{i,k-2}, a_{i,k-2}^{\text{CL}}, a_{i-1,k-2}^{\text{CL}}), a_{i,k-1}^{\text{CL}}, a_{i-1,k-1}^{\text{CL}}) \\
& \dots \\
& = f(x_{i,k-\tau_i,k}, \{a_{i,k'}^{\text{CL}}\}_{k'=k-\tau_i,k}^{k-1}, \tau_{i,k}, \{a_{i-1,k'}^{\text{CL}}\}_{k'=k-\tau_i,k}^{k-1}) \\
& = f(S_{i,k}^{\text{CL}}, \{a_{i,k'}^{\text{CL}}\}_{k'=k-\tau_i,k}^{k-1})
\end{aligned} \tag{43}$$

Therefore, if the sequence of actions of vehicle $i - 1$ within the time window of $[k - \tau_{i,k} + 1, k]$ is available to vehicle i , the current driving status $x_{i,k}$ can be fully determined by the delayed observation $S_{i,k}^{\text{CL}}$.

We would like to point out that the state space model in (41) ignores the random perturbations caused by mechanical noise due to practical system limitations. Therefore, the practical driving status $x_{i,k}$ might slightly deviate from the derived values by (43). However, the sequence of PC actions $\{a_{i,k'}^{\text{CL}}\}_{k'=k-\tau_{i,k}}^{k-1}$ of vehicle $i - 1$ still play the most important role in predicting the current driving status $x_{i,k}$ based on the delayed observation $S_{i,k}^{\text{CL}}$.

E. Proof of Theorem 2

From Lemma 2 and Lemma 3, we can derive that if vehicle $i - 1$'s control actions $a_{i-1,k'}^{\text{CL}}$ within the period $k' \in [k - \tau_{\max}, k - 1]$ is available to vehicle i at control interval k , the optimal PC policy $\pi_i^{\text{CL}*}(S_{i,k}^{\text{CL}})$ based on delayed observation $S_{i,k}^{\text{CL}}$ performs as well as the optimal policy $\hat{\pi}_i^{\text{CL}*}(x_{i,k})$ based on the current driving status $x_{i,k}$ at vehicle i . In practice, only the control action $a_{i-1,k-\tau_{i,k}}^{\text{CL}}$ of vehicle $i - 1$ is available to vehicle i at control interval k . Consider the extreme case when $a_{i-1,k'}^{\text{CL}}$ is a constant value within period $k' \in [k - \tau_{\max}, k - 1]$, which means that $a_{i-1,k-\tau_{i,k}}^{\text{CL}}$ can accurately represent all of vehicle $i - 1$'s control actions within the time window. In this case, $x_{i,k}$ can be fully determined by $S_{i,k}^{\text{CL}}$ according to Lemma 3. Thus, $\hat{\pi}_i^{\text{CL}*}(x_{i,k})$ and $\pi_i^{\text{CL}*}(S_{i,k}^{\text{CL}})$ are equivalent from Lemma 2. In other words, if the variation of $a_{i-1,k'}^{\text{CL}}$ is zero within period $k' \in [k - \tau_{\max}, k - 1]$, the PC performance at vehicle i is not impacted by the observation delay. Extending to the general case, the smaller the variation of $a_{i-1,k'}^{\text{CL}}$ within the period $k' \in [k - \tau_{\max}, k - 1]$, the more accurate $a_{i-1,k-\tau_{i,k}}^{\text{CL}}$ can be used to represent all the control actions of vehicle $i - 1$ within the time window, and the smaller the impact of observation delay on PC performance of vehicle i .

REFERENCES

- [1] S. E. Li, Y. Zheng, K. Li, Y. Wu, J. K. Hedrick, F. Gao, and H. Zhang, "Dynamical modeling and distributed control of connected and automated vehicles: Challenges and opportunities," *IEEE Intelligent Transportation Systems Magazine*, vol. 9, no. 3, pp. 46–58, 2017.
- [2] L. Lei, T. Liu, K. Zheng, and L. Hanzo, "Deep reinforcement learning aided platoon control relying on v2x information," *IEEE Transactions on Vehicular Technology*, vol. 71, no. 6, pp. 5811–5826, 2022.
- [3] L. Lei, Y. Tan, K. Zheng, S. Liu, K. Zhang, and X. Shen, "Deep reinforcement learning for autonomous internet of things: Model, applications and challenges," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 1722–1760, 2020.
- [4] S. Omidshafiei, J. Papis, C. Amato, J. P. How, and J. Vian, "Deep decentralized multi-task multi-agent reinforcement learning under partial observability," in *Proc. International Conference on Machine Learning*, PMLR, 2017, pp. 2681–2690.
- [5] H. Ye, G. Y. Li, and B.-H. F. Juang, "Deep reinforcement learning based resource allocation for v2v communications," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 4, pp. 3163–3173, 2019.
- [6] H. Yang, X. Xie, and M. Kadoch, "Intelligent resource management based on reinforcement learning for ultra-reliable and low-latency iov communication networks," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 5, pp. 4157–4169, 2019.
- [7] K. Zia, N. Javed, M. N. Sial, S. Ahmed, A. A. Pirzada, and F. Pervez, "A distributed multi-agent rl-based autonomous spectrum allocation scheme in d2d enabled multi-tier hetnets," *IEEE Access*, vol. 7, pp. 6733–6745, 2019.
- [8] Z. Nan, Y. Jia, Z. Ren, Z. Chen, and L. Liang, "Delay-aware content delivery with deep reinforcement learning in internet of vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 7, pp. 8918–8929, 2021.
- [9] 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.
- [10] H. V. Vu, Z. Liu, D. H. Nguyen, R. Morawski, and T. Le-Ngoc, "Multi-agent reinforcement learning for joint channel assignment and power allocation in platoon-based c-v2x systems," *arXiv preprint arXiv:2011.04555*, 2020.
- [11] P. Xiang, H. Shan, M. Wang, Z. Xiang, and Z. Zhu, "Multi-agent rl enables decentralized spectrum access in vehicular networks," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 10, pp. 10750–10762, 2021.
- [12] H. Zhang, C. Lu, H. Tang, X. Wei, L. Liang, L. Cheng, W. Ding, and Z. Han, "Mean-field aided multi-agent reinforcement learning for resource allocation in vehicular networks," *IEEE Internet of Things Journal*, 2022.
- [13] 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.
- [14] Y. Xu, K. Zhu, H. Xu, and J. Ji, "Deep reinforcement learning for multi-objective resource allocation in multi-platoon cooperative vehicular networks," *IEEE Transactions on Wireless Communications*, 2023.
- [15] Y. Zhong, B. Wang, and Y. Wang, "Reward backpropagation prioritized experience replay," 2017.
- [16] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [17] T. Liu, L. Lei, K. Zheng, and K. Zhang, "Autonomous platoon control with integrated deep reinforcement learning and dynamic programming," *IEEE Internet of Things Journal*, vol. 10, no. 6, pp. 5476–5489, 2023.
- [18] Y. S. Nasir and D. Guo, "Deep reinforcement learning for distributed dynamic power allocation in wireless networks," *arXiv preprint arXiv:1808.00490*, vol. 8, p. 2018, 2018.
- [19] S. Kaul, R. Yates, and M. Gruteser, "Real-time status: How often should one update?" in *Proc. IEEE International Conference on Computer Communications (INFOCOM)*. IEEE, 2012, pp. 2731–2735.
- [20] A. Kosta, N. Pappas, A. Ephremides, and V. Angelakis, "Age and value of information: Non-linear age case," in *Proc. IEEE International Symposium on Information Theory (ISIT)*. IEEE, 2017, pp. 326–330.
- [21] Z. Wang, M.-A. Badiu, and J. P. Coon, "A value of information framework for latent variable models," in *Proc. IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2020, pp. 1–6.
- [22] Z. Wang, M.-A. Badiu, and J. P. Coon, "A framework for characterizing the value of information in hidden markov models," *IEEE Transactions on Information Theory*, vol. 68, no. 8, pp. 5203–5216, 2022.
- [23] O. Ayan, M. Vilgelm, M. Klügel, S. Hirche, and W. Kellerer, "Age-of-information vs. value-of-information scheduling for cellular networked control systems," in *Proc. 10th ACM/IEEE International Conference on Cyber-Physical Systems*, 2019, pp. 109–117.
- [24] J. Foerster, N. Nardelli, G. Farquhar, T. Afouras, P. H. Torr, P. Kohli, and S. Whiteson, "Stabilising experience replay for deep multi-agent reinforcement learning," in *Proc. International conference on machine learning*. PMLR, 2017, pp. 1146–1155.
- [25] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *Proc. The AAAI conference on artificial intelligence*, vol. 30, no. 1, 2016.
- [26] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski et al., "Human-level control through deep reinforcement learning," *nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [27] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," *Advances in neural information processing systems*, vol. 30, 2017.