# Multi-Agent Deep Reinforcement Learning in Vehicular OCC

Amirul Islam, Leila Musavian, Nikolaos Thomos
CSEE, University of Essex, UK.
Email: {amirul.islam, leila.musavian, nthomos}@essex.ac.uk

*Abstract*—Optical camera communications (OCC) has emerged as a key enabling technology for the seamless operation of future autonomous vehicles. In this paper, we introduce a spectral efficiency optimization approach in vehicular OCC. Specifically, we aim at optimally adapting the modulation order and the relative speed while respecting bit error rate and latency constraints. As the optimization problem is NP-hard problem, we model the optimization problem as a Markov decision process (MDP) to enable the use of solutions that can be applied online. We then relaxed the constrained problem by employing Lagrange relaxation approach before solving it by multi-agent deep reinforcement learning (DRL). We verify the performance of our proposed scheme through extensive simulations and compare it with various variants of our approach and a random method. The evaluation shows that our system achieves significantly higher sum spectral efficiency compared to schemes under comparison.

*Index Terms*—Deep reinforcement learning, optical camera communication, vehicular communication, Lagrangian relaxation

## I. INTRODUCTION

Autonomous vehicles are driving the revolution in future smart cities and are considered as the leading transformative technologies in intelligent transportation systems (ITS). To cope with the current ever-growing and complex nature of vehicular networks, data sharing on the road involves continuously increasing amounts of data and thus incurring enormous network overhead [1]. This puts tremendous pressure on the overused radio frequency (RF) spectrum that is already congested and saturated. On the contrary, the recent advancements and potential advantages of optical camera communication (OCC) over RF-based communication systems, such as license-free unlimited spectrum, lower implementation cost, and enhanced security, have rendered this technology to be an essential alternative for ITS [2], [3]. OCC uses light-emitting diodes (LEDs) as transmitters and cameras as receivers.

One of the main challenges of vehicular networks is that they are highly dynamic and require processing of huge amounts of data. The effectiveness of the ITS depends on supporting vehicle-to-vehicle (V2V) communication within the shortest time and lowest error. Recently, several technologies have been explored for ITS, which target delay minimization [4], reliability maximization [5] using traditional distributed method to solve the underlying optimization problems. In [4], the transmission power of the vehicular network is minimized by grouping vehicles into clusters and defining reliability as queuing delay violation probability. In

[5], a joint resource allocation and power control algorithm is proposed to maximize the communication rate considering latency and reliability constraints. However, meeting the required reliability target and at the same time respecting stringent time constraints makes the V2V communication more challenging. In particular, these problems are difficult to solve following the traditional distributed methods because of their complexity and the entailed time needed. Notably when it involves decision-making in controlling different parameters, e.g., speed, distance and modulation schemes.

Reinforcement learning (RL) can serve as an effective alternative solution to overcome the complexity of such problems [6]. Specifically, RL offers decentralized decision-making when the centralized decision is impossible to make. In this paper, we adopt RL, and hence, we first model the studied problem as a Markov Decision Process (MDP). Methods like value iteration that is commonly proposed to solve the MDP, require knowing the state transition probabilities beforehand making it difficult to evaluate the optimal policy. These complexities are overcome through using Q-Learning [6]. However, Q-Learning is characterized by slow convergence rate, and hence, is inappropriate for solving large-scale problems as the ones we study here. To address this limitation of the Q-Learning algorithm, we use deep RL (DRL) [7].

Despite overcoming the problem of Q-Learning, DRL faces difficulties in solving large-scale V2V networks in a centralized way. This introduces higher latency, which may result in increased failure rates as the vehicles may make a decision using outdated information, which eventually compromises safety. To address this issue, we utilize the concept of independent learning and multi-agent RL (MARL) [8]. In independent-learning MARL, each agent learns its policy independently and exploiting by a local observation while modelling other agents as parts of the environment dynamics.

To the best of our knowledge, DRL-based performance optimization in vehicular OCC has not been investigated in the literature. Several studies suggest the use of RL in hybrid RF and photodiode (PD)-based visible light communications (VLC) networks [9], [10]. The authors in [9] apply reinforcement learning for network selection by considering the traffic type and the possibility of having learning records to improve the Q-Learning algorithm. In [10], the authors implement MARL to develop online power allocation. These works are interesting; however, they consider a centralized DRL scheme and ignore the inherent latency and reliability requirements.

Moreover, they study PD-based receiver, which faces interference problems when dealing with multiple vehicles. OCC overcomes interference problems as it can spatially separate and process different transmitter sources independently on its image plane, which has millions of pixels, and this provides freedom to handle multiple users.

In this paper, we propose independent and multi-agent DRL-based spectral efficiency maximization scheme in vehicular OCC. We maximize the spectral efficiency by adapting the modulation order from a chosen set of available modulation schemes and also by changing the relative speed of the agent (vehicle) while satisfying bit error rate (BER) and latency constraints. To the best of our knowledge, we introduce DRL for the first time in vehicular OCC for optimizing the spectral efficiency. The major contributions of this paper are summarized as follows:

- We formulate spectral efficiency maximization problem subject to BER, latency and a small set of modulation orders in vehicular OCC. The optimization function is a non-deterministic polynomial-time (NP) hard problem leading to a difficult search for the optimal solution. Hence, we first model the problem as an MDP;
- We relax the constrained maximization problem by converting it to an unconstrained problem using the Lagrangian relaxation method and then solve it using deep Q-Learning;
- We evaluate the proposed DRL-based optimization scheme and compare it with variants of the proposed scheme and a random scheme. The results show that it can effectively learn how to maximize the spectral efficiency while meeting the constraints and our solution outperforms significantly the other schemes.

## II. SYSTEM MODEL AND PROBLEM FORMULATION

### A. OCC System Model

We consider a vehicular OCC system model as shown in Fig. 1, where each vehicle is an individual agent. Each vehicle has a transmitting unit at the back consisting of LEDs backlights and a vision camera set and a receiving unit at the front having a high-speed camera (1000 frame per second (fps)). The camera at the back measures the backward distance using a stereo-vision camera. We consider $B$ to be the number of V2V links in the back of each vehicle and $\mathcal{B} = \{1, 2, \cdots B\}$ the set of V2V links. We express the distance of the agent (vehicle) with the backward vehicles as $d^b$, where $b \in \mathcal{B}$ is the index of the V2V link.[1]

We introduce an adaptive modulation scheme of M-ary quadrature amplitude modulation (M-QAM) as it can offer low BER and improved spectral efficiency [11]. Similar to [12], we use time division multiple access (TDMA) in our system to transmit at different modulation orders for different backward vehicles. In TDMA, each link of the vehicle transmits at a specific time only. Hence, the spectral efficiency is divided by the number of available users, $B$, at the back.

[1]In this paper, $\mathcal{B}$ denotes the set of V2V links.
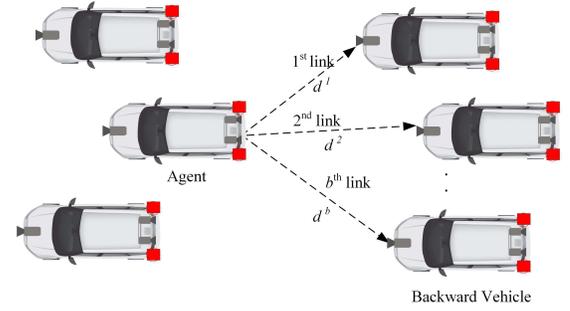


Fig. 1: Proposed system model for vehicular optical camera communication.

To ensure the vehicles are free from obstruction and that can continuously communicate with each other, we assume that the transmitter and camera receiver has an uninterrupted line-of-sight link between them. The channel gain, $H_t^b$, of link $b$ in time $t$ is expressed similarly to [13] as:

$$H_t^b = \begin{cases} \frac{(m+1)A}{2\pi(d_t^b)^2} \cos^m(\phi) \, T_s(\theta) \, g \, \cos(\theta), & 0 \leq \theta \leq \theta_l \\ 0, & \theta > \theta_l \end{cases} \quad (1)$$

where $m$ is the the order of the Lambertian radiation pattern, which is derived from LED semi-angle at half luminance, $\Phi_{1/2}$, as $m = \frac{-\ln(2)}{\ln(\cos(\Phi_{1/2}))}$. $A$ is area of the entrance pupil of the camera lens, $d_t^b$ is the agents' distance with the backward vehicles at time $t$, $\phi$ is the angle of irradiance with respect to the emitter, $T_s(\theta)$ corresponds to the transmission efficiency of the optical filter, $g$ is the gain of the lens, $\theta$ is the angle of incidence (AoI) with respect to the receiver axis, and $\theta_l$ denotes the FoV of the image sensor lens. An ideal lens has a gain: $g = n^2/\sin^2(\theta_l)$, where $n$ is the internal refractive index of the lens.

### B. Performance Parameter Definition

In this subsection, we specify the performance defining metrics of OCC in terms of signal-to-noise ratio (SNR), the achievable rate, and the observed transmission latency. First, we express SNR to define the communication link quality of the signal transmission. In particular, the received SNR, $\gamma_t^b$, of the link $b$ in time $t$ for a single LED-camera communication is expressed similarly to [14] as:

$$\gamma_t^b = \frac{\left(\rho P_{r,t}^b\right)^2}{\sigma_t^b} = \frac{\left(\rho H_t^b P\right)^2}{\sigma_t^b} \quad (2)$$

where $\rho$ is the receiver's responsivity, $P$ is the optical transmit power, and $\sigma_t^b$ represents the total noise power, which is written as:

$$\sigma_t^b = q\rho P_n A_t^b W_{\text{fps}} , \quad (3)$$

where $q$ is electron charge, $P_n$ is the background noise power per unit area, $A_t^b$ area of the receiver for the link $b$ at time $t$, and $W_{\text{fps}}$ is the sampling rate of the camera in fps. We can calculate $A_t^b$ using the similar concept in [15]. To remove the effect of quantization in the received signal, measurements are made at the square grid of points. Therefore, the LED will

occupy a square area of size $l'^2 = A$ having the diameter of $l' = fl/d^b$ of a circle since LEDs will form a circular shape on the receiver grid, where $l$ is the diameter of a LED, and $f$ is the focal length. When the projected diameter of the image becomes smaller than the size of a pixel, we refer to this as critical distance, $d_c = fl/s$, where $s$ is the edge-length of a pixel. Based on the above definitions and from (3), (2) can be summarized as,[2]

$$\gamma^b = \begin{cases} \frac{\rho k^2 P^2}{q P_n W f^2 l^2 (d^b)^2}; & \text{if } d^b < d_c , \\ \frac{\rho k^2 P^2}{q P_n W s^2 (d^b)^4}; & \text{if } d^b \geq d_c . \end{cases} \quad (4)$$

where $k = \frac{(m+1)A}{2\pi} \cos^m(\phi) \, T_s(\theta) \, g \, \cos(\theta)$.

We evaluate the BER of the optical wireless channel at the receiver using the M-QAM scheme similar to [16]. Considering M-QAM, the spectral efficiency is expressed as, $\text{SE}^b = \log_2(M^b)$, where $M^b$ is the available constellation points for each V2V link, $b$, e.g., $M = 4, 8, 16, \cdots$.

The channel capacity (measured in bits/sec) of the camera-based communication system is derived from the employed modulation scheme of link $b$ as has been shown in [17] as

$$C^b = \frac{(W_{\text{fps}}/3) N_{\text{LEDs}} w \varrho}{2 \tan\left(\frac{\theta_l}{2}\right) \cdot d^b} \cdot \log_2(M^b), \quad (5)$$

where $N_{\text{LEDs}}$ is the number of LEDs at each row of the transmitter, $w$ is the image width, and $\varrho$ is the size of LED lights in cm$^2$. Please note that, the distance $d^b$ in (5) is affected by relative speed of the vehicle $v$, which also affects the position of the vehicle. The inter-vehicular distance at current time $t$ is adjusted using $d_t = d_{t-1} + v_t \cdot \Delta t$, where $d_{t-1}$ is the distance at previous time instant and $\Delta t$ is the time elapsed between time instants $t$ and $t-1$.

We consider that the end-to-end latency is dominated by the transmission latency, and we neglect the computational latency. This is because we process a small amount of data, and hence, the computational time will be short. Thus, the transmission latency of packet size, $L$, is expressed as $\tau^b = L/C^b$.

### C. Problem Formulation

Considering the proposed framework, we formulate a sum spectral efficiency optimization scheme. We aim at selecting modulation scheme from the available set of modulation orders and controlling the relative speed of the vehicle. The BER and latency are constrained so that they meet the values imposed by the system. Finally, the proposed maximization problem is formulated as:

$$\max_{\mathcal{M}, v} \quad \frac{1}{B} \sum_{b=1}^{B} \log_2\left(M^b\right), \quad (6)$$

$$\text{s.t.} \quad \text{BER}^b \leq \text{BER}_{\text{tgt}}, \; \forall b; \quad (7)$$

$$\tau^b \leq \tau_{\max}, \; \forall b; \quad (8)$$

$$M^b \in \mathcal{M}, \; \forall b; \quad (9)$$

[2]For notational simplicity, we drop $t$ from the notation in the remainder of the paper unless it is necessary; hence, we will adopt $\gamma^b$ instead of $\gamma_t^b$ and so on. Also, it is clear from the context that distance is our working variable.

where $\mathcal{M}$ is the set of the available modulation orders, $\text{BER}_{\text{tgt}}$ is the maximum target BER, and $\tau_{\max}$ is the maximum allowed latency. Equations (7) and (8) correspond to the BER and latency constraints. The modulation scheme is chosen from a small set of available M-QAM schemes, as shown in (9).

### III. DRL-BASED PROBLEM FORMULATION AND PROPOSED SOLUTION

The optimization problem presented in (6) - (9) is an NP-hard combinatorial problem [18], and thus it is challenging to find the optimal solution. The optimization problem also includes non-linear operations, such as, (6) and (8). To solve this problem using a distributed method, each agent should choose the speed and modulation scheme separately, which makes the solution more complex and time-consuming. Recall that, in vehicular communication, we should meet low latency and required BER targets to ensure that the information is received reliably within the shortest time. Reinforcement learning is an effective way to solve such dynamic and time-varying problems as it can learn the optimal policy through interaction with the environment, and this way adjust to the environmental changes over time.

### A. Modeling of MDP

The optimization problem (6) is modelled as an MDP with a tuple $(\mathcal{S}, \mathcal{A}, p, r, \zeta)$ [6], where $\mathcal{S}$ is the set of all possible states; $\mathcal{A}$ denotes the set of all possible actions; $p$ denotes the transition probability $p(s_{t+1}, r_t | s_t, a_t)$ when the agent selects an action $a_t \in \mathcal{A}$ and transits to a new state $s_{t+1} \in \mathcal{S}$ from the current state $s_t \in \mathcal{S}$; and $r$ represents the reward. While $\zeta \in [0, 1]$ is the discount factor, which gradually discounts the effect of an action to future rewards. The state space $\mathcal{S}$, the action space $\mathcal{A}$, and the reward function, $r$ of the considered RL framework are defined below.

**State**: The observed states of our considered environment include: the backward distance vector, $\mathbf{d}_t^b = (d_t^1, \cdots, d_t^B)$ and the modulation order set, $\mathcal{M} = \{4, 8, 16, 32, 64\}$.

**Action:** At state $s_t$, the agent takes an action $a_t$, by changing its relative speed, $v_t$ and selecting the modulation order from the set $\mathcal{M}$.

**Reward:** The reward function that guides the overall learning should be consistent with the objective. Since our objective is to maximize the sum spectral efficiency, we design our reward function as a weighted sum of a reward related to the backward distance and the sum spectral efficiency (6). First, we express the reward related to distance as follows:

$$r_t^{\text{d},i} = \begin{cases} -1 \times (d_{\text{stop}} - d_t^b), & d_t^b < d_{\text{stop}} , \\ \frac{1}{d_t^b - d_{\text{stop}}}, & d_t^b > d_{\text{stop}} , \end{cases} \quad (10)$$

where $i$ is the index of the agent and $d_{\text{stop}}$ is the stopping distance [19]. In our system, each vehicle will carry out the same process individually. As a result, for notational simplicity, we drop $i$ hereafter.

Finally, considering the objective function of (6), the overall reward, $R_t$, can be expressed as

$$R_t = \omega_d \, r_t^{\mathrm{d}} + \omega_r \, \frac{1}{B} \sum_{b=1}^{B} \log_2 \left( M_t^b \right), \qquad (11)$$

where $\omega_d$ and $\omega_r$ are positive weights that balance distance and sum spectral efficiency rewards. The weights can be adjusted based on the system requirements.

### B. RL-based Problem Formulation

The goal of RL is to maximize the expected return from the state $s_t$ by determining the optimal policy. The return, $G_t$, is defined as the cumulative discounted reward, and is expressed as follows:

$$G_t = \sum_{j=0}^{\infty} \zeta^j R_{t+j+1}, \qquad 0 \le \zeta \le 1. \qquad (12)$$

In summary, the objective of our proposed system is to determine the optimal policy, i.e., to select the speed and modulation order while respecting the BER and latency constraints. From the above consideration, the reward maximization problem that corresponds to the problem formulation presented in Section II-C is expressed as

$$\max \quad \mathbb{E}\left[ G_t \left( s_t, a_t \right) \right], \; \forall t \qquad (13)$$

$$\text{s.t.} \quad \mathrm{BER}_t^b \le \mathrm{BER}_{\mathrm{tgt}}, \; \forall t; \qquad (14)$$

$$\tau_t^b \le \tau_{\max}, \; \forall t; \qquad (15)$$

### C. Solution of the Problem

We can solve the constrained MDP problem presented in (13)-(15) by converting it into an unconstrained one following Lagrange Relaxation method [20]. Therefore, by relaxing the BER and latency constraints, we re-express the constrained optimization problem as:

$$c^{\boldsymbol{\lambda},\boldsymbol{\nu}}(s_t, a_t) = R_t \left( s_t, a_t \right) - \sum_{b=1}^{B} \lambda^b \cdot (\mathrm{BER}_t^b - \mathrm{BER}_{\mathrm{tgt}})$$

$$- \sum_{b=1}^{B} \nu^b \cdot (\tau_t^b - \tau_{\max}), \qquad (16)$$

where $\boldsymbol{\lambda} = (\lambda^1, \lambda^2 \cdots, \lambda^b)$ and $\boldsymbol{\nu} = (\nu^1, \nu^2, \cdots, \nu^b)$ are vectors representing the Lagrange multipliers corresponding to the constraints in (14) and (15), respectively. The optimal value of the constrained MDP problem is computed as [21]:

$$L_\delta^{\pi^*,\boldsymbol{\lambda}^*,\boldsymbol{\nu}^*}(s) = \max_{\boldsymbol{\lambda},\boldsymbol{\nu} \ge 0} \min_{\pi \in \phi} V^{\pi,\boldsymbol{\lambda},\boldsymbol{\nu}}(s) - \sum_{b=1}^{B} \lambda^b \delta_1 - \sum_{b=1}^{B} \nu^b \delta_2, \qquad (17)$$

where $\delta = \{\delta_1, \delta_2\}$, with $\delta_1 = \mathrm{BER}_{\mathrm{tgt}}$ and $\delta_2 = \tau_{\max}$. $\phi$ denotes the set of all possible stationary policies.

In practice, the optimal policy cannot be determined using value iteration as it requires knowledge of transition probabilities beforehand, which is not possible because of the size

TABLE I: Vehicular OCC modelling parameters

| Parameter, Notation | Value |
| --- | --- |
| Angle of irradiance w.r.t. the emitter, $\phi$ | $70^o$ |
| AoI w.r.t. the receiver axis, $\theta$ | $60^o$ |
| FOV of the camera lens, $\theta_l$ | $90^o$ |
| Image sensor physical area, $A$ | $10 \text{ cm}^2$ |
| Optical filter Transmission efficiency, $T_s$ | 1 |
| Concentrator/lens gain, $g$ | 3 |
| Optical transmitting power, $P$ | 1.2 Watts |
| Modulation scheme set, $\mathcal{M}$ | 4, 8, 16, 32, 64 |
| Camera-frame rate, $W_{\mathrm{fps}}$ | 1000 fps |
| Number of LEDs at each row, $N_{\mathrm{LEDs}}$ | 30 |
| Packet size, $L$ | 5 kbits |
| Size of the LED, $\varrho$ | $15.5 \times 5.5 \text{ cm}^2$ |
| Resolution of image, $w$ | $512 \times 512$ pixels |

of the state and action space.[3] To solve this problem, we adopt a model-free RL approach known as tabular Q-Learning. This approach uses the $Q_t(s_t, a_t)$ values for each state-action pair instead of the value function, which is a function of the state only. The Q-Learning algorithm employs the following recursive formula to update the $Q_t(s_t, a_t)$ values:

$$Q_{t+1}(s_t, a_t) = (1 - \alpha_t) Q_t(s_t, a_t) + \alpha_t \left[ c_t(s_t, a_t) + \zeta \max_{a_{t+1} \in \mathcal{A}} Q_t(s_{t+1}, a_{t+1}) \right], \qquad (18)$$

where $\alpha_t \in [0, 1]$ is a time-varying learning rate. After computing the Q-values, the optimal policy $\pi^*$ can be determined as

$$\pi^{*,\boldsymbol{\lambda},\boldsymbol{\nu}}(s_t) = \arg \max_{a_t \in \mathcal{A}} Q^{*,\boldsymbol{\lambda},\boldsymbol{\nu}}(s_t, a_t), \forall s \in S. \qquad (19)$$

When the state-action space is small, Q-Learning can determine the optimal policy. However, Q-Learning cannot decide the value functions or optimal policies accurately in large-scale problems within a reasonable time because of the state-action space. This problem can be solved by employing deep learning-based function approximators by means of deep neural networks. Deep Q-network (DQN) is used to train the network and learn the optimal policy.

In order to stabilize the learning of DQN, we follow the target network approach. The DQN consists of two separate networks known as the main network that approximates the Q-function and the target network that gives the target for updating the main network. The target network is not updated after each iteration because it adjusts the main network updates to control the value estimations. If both networks are updated simultaneously, the change in the main network would be exaggerated due to the feedback loop from the target network, which results in an unstable network.

To ensure convergence, the neural network aims to minimize the loss function, $L(\boldsymbol{\beta})$, which can be defined as

$$L(\boldsymbol{\beta}) = \mathbb{E}\left[ y_t - Q\left( s_t, a_t; \boldsymbol{\beta} \right) \right]^2, \qquad (20)$$

---

[3]For notational simplicity, we drop the Lagrangian multipliers from the notation in the remainder of the paper unless it is necessary, for example, we will write $c(s_t, a_t)$, $Q^*(s_t)$, instead of $c^{\boldsymbol{\lambda},\boldsymbol{\nu}}(s_t, a_t)$, $Q^{*,\boldsymbol{\lambda},\boldsymbol{\nu}}(s_t)$, respectively.

TABLE II: List of DRL hyper-parameters and their values

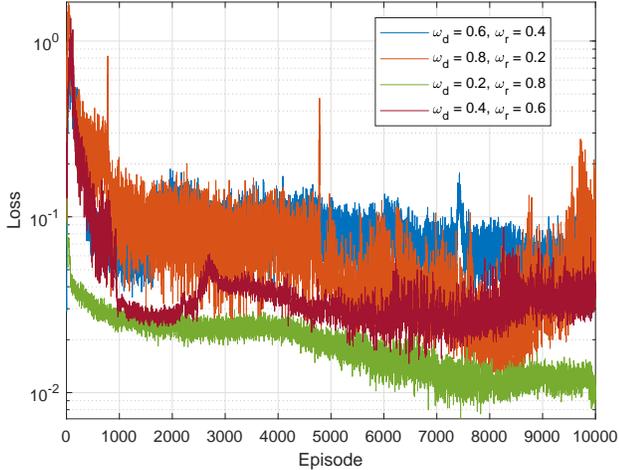| Parameter, Notation | Value |
|---|---|
| Mini-batch size | 32 |
| Replay memory size | 100000 |
| Number of hidden layer (Neurons) | 1(250) |
| Exploration rate, $\epsilon$ | 0.05 |
| Discount factor, $\zeta$ | 0.98 |
| Activation function | ReLU |
| Optimizer | RMSProp |
| Learning rate (used by RMSProp) | 0.001 |
| Gradient momentum (used by RMSProp) | 0.95 |



Fig. 2: Convergence of loss function for $\epsilon = 0.05$ and learning rate $\alpha = 0.001$.

where $y_t = c(s_t, a_t) + \zeta \max_{a_{t+1} \in \mathcal{A}} Q(s_{t+1}, a_{t+1}; \boldsymbol{\beta}_-)$ is the target for each iteration. $\boldsymbol{\beta}$ denotes the neural network's parameters of current iteration and $\boldsymbol{\beta}_-$ is the value from the previous update. Note that, $\boldsymbol{\beta}_-$ are held fixed when optimizing the loss function $L(\boldsymbol{\beta})$. The optimal value of the Lagrange multipliers $\lambda^b$, $\nu^b$ in (16) can be learned online using a stochastic sub-gradient method as presented in [22].

## IV. SIMULATION SETUP AND RESULTS

### A. Simulation Setup

To evaluate the performance of the proposed system, we build a simulation environment upon traffic simulator Simulation of Urban Mobility (SUMO) [23]. Our simulation framework maintains the connection between SUMO and the DRL agent using Traffic Control Interface (TraCI). The DQN consists of three fully connected layers, including an input layer, a hidden layer, and an output layer. The hidden layer has 250 neurons. We use rectified linear unit (ReLU) as the activation function. We then adopt root mean square propagation (RMSProp) optimizer to minimize the loss, where we set an initial learning rate to 0.001, which varies over time. The neural network is designed using Tensorflow [24]. We implement $\epsilon$-greedy policy to balance between exploration and exploitation.

In our simulation, we train the DQN for 10000 episodes. The exploration rate, $\epsilon$, is set to 0.05. We choose a discount factor, $\zeta = 0.98$. We present the simulation parameters for
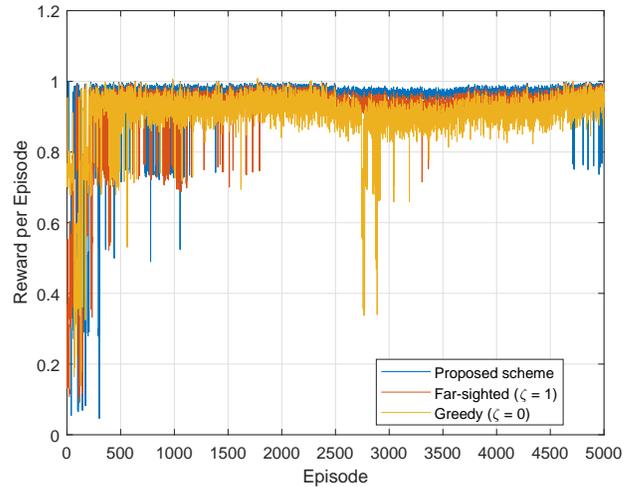


Fig. 3: Reward per training episode for three different approaches for $\epsilon = 0.05$ and learning rate $\alpha = 0.001$.

the OCC system model in Table I, whereas the training and testing parameters of the DRL are listed in Table II.

We investigate the performance of the proposed multi-agent DRL-based vehicular scheme against different methods for comparison, namely, greedy, far-sighted and random scheme. In greedy method, we assume $\zeta = 0$ in (20), whereas in far-sighted case, we assume $\zeta = 1$, while we keep all other parameters of the systems as reported in Table II. Finally, in the random scheme, the action is chosen randomly for all the vehicles at each time slot.

### B. Simulation Results

First, we perform an ablation study to determine the optimal weight values for distance and spectral efficiency rewards in (11). In doing so, we examine the proposed algorithm for different weight settings, but for simplified representation, we include four settings, such as $\omega_d = 0.2$ and $\omega_r = 0.8$, $\omega_d = 0.4$ and $\omega_r = 0.6$, $\omega_d = 0.6$ and $\omega_r = 0.4$, $\omega_d = 0.8$ and $\omega_r = 0.2$, as shown in Fig. 2. We can see that we achieve improved loss performance when we allocate more weight value toward the spectral efficiency part. The figure points that the algorithm converges at 8000 episodes for $\omega_d = 0.2$ and $\omega_r = 0.8$. On the contrary, other weight needs more time to converge and show frequent variations in the loss and offer higher loss than $\omega_d = 0.2$ and $\omega_r = 0.8$ set. Thus, we adopt these weight values for the rest of our performance evaluation.

To analyze the convergence behaviour of the multi-agent vehicular OCC system, we demonstrate the cumulative rewards per episode for three different discount factors, i.e., the proposed scheme ($\zeta = 0.98$), greedy ($\zeta = 0$) and far-sighted ($\zeta = 1$). The results are shown in Fig. 3. From this figure, we observe that until 1500 episodes, the greedy and far-sighted approaches perform better than the proposed scheme. However, the cumulative reward for the proposed schemes improves as the training advances and eventually reaches to lower loss. We can conclude that the proposed scheme achieves higher rewards than all the schemes under
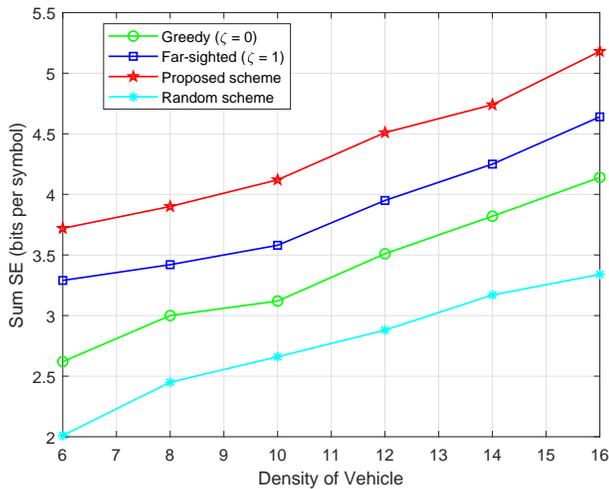
Fig. 4: Comparison of sum spectral efficiency with different approaches for $\epsilon = 0.05$ and learning rate $\alpha = 0.001$.

consideration.

Finally, we present the maximized sum spectral efficiency versus various density of vehicles for all schemes under comparison in Fig. 4. From the figure, we see that the sum spectral efficiency increases with an increase in the density of vehicles for all the methods, and there is a significant performance gap between each scheme. We also observe that our proposed scheme can achieve a maximum of 5.2 bits per symbol, whereas the random method can achieve 3.3 bits per symbol. The results show that the proposed algorithm obtains approximately 2.3 times better rates in comparison to the random scheme, 1.25 times for far-sighted, and about 1.11 times for greedy schemes when the density of vehicles is 16. Accordingly, we can conclude that our proposed OCC system outperforms all the other schemes.

## V. CONCLUSION

In this paper, we present a multi-agent DRL-based spectral efficiency optimization scheme in vehicular OCC while respecting BER and latency requirements. In doing so, we optimize our system by selecting the optimal modulation order and adjusting the relative speed of each vehicle. To overcome the inherent complexity of the studied problem, we model the problem as an MDP. We then convert the constrained problem into an unconstrained problem using the Lagrangian relaxation method. Next, we solve the problem by employing deep Q-Learning to deal with the large state-action spaces we encounter. Finally, we verify the performance of our scheme through extensive simulations and compare it with various variants of our scheme. The evaluations reveal that our system achieves better sum spectral efficiency compared to the schemes under comparison.

## REFERENCES

[1] P. Papadimitratos, A. De La Fortelle, K. Evenssen, R. Brignolo, and S. Cosenza, "Vehicular communication systems: Enabling technologies, applications, and future outlook on intelligent transportation," *IEEE Commun. Mag.*, vol. 47, no. 11, pp. 84–95, Nov. 2009.

[2] I. Takai, T. Harada, M. Andoh, K. Yasutomi, K. Kagawa, and S. Kawahito, "Optical vehicle-to-vehicle communication system using LED transmitter and camera receiver," *IEEE Photon. J.*, vol. 6, no. 5, pp. 1–14, Oct. 2014.

[3] T. Y. ín *et al.*, "Image-sensor-based visible light communication for automotive applications," *IEEE Commun. Mag.*, vol. 52, no. 7, pp. 88–97, Jul. 2014.

[4] M. I. Ashraf, C.-F. Liu, M. Bennis, and W. Saad, "Towards low-latency and ultra-reliable vehicle-to-vehicle communication," in *Proc. EuCNC'17*, Oulu, Finland, Jun. 2017, pp. 1–5.

[5] W. Sun, E. G. Ström, F. Brännström, Y. Sui, and K. C. Sou, "D2D-based V2V communications with latency and reliability constraints," in *Proc. 2014 IEEE GC Wkshps*, Austin, TX, USA, Dec. 2014, pp. 1414–1419.

[6] R. S. Sutton and A. G. Barto, *Introduction to reinforcement learning*. MA, USA: MIT press Cambridge, 1998, vol. 135.

[7] H. Li, T. Wei, A. Ren, Q. Zhu, and Y. Wang, "Deep reinforcement learning: Framework, applications, and embedded implementations," in *Proc. ICCAD'17*, Irvine, CA, USA, Nov. 2017, pp. 847–854.

[8] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in *Proc. ICML'93*, CA, USA, Jul. 1993, pp. 330–337.

[9] Z. Du, C. Wang, Y. Sun, and G. Wu, "Context-aware indoor VLC/RF heterogeneous network selection: Reinforcement learning with knowledge transfer," *IEEE Access*, vol. 6, pp. 33 275–33 284, Jun. 2018.

[10] J. Kong, Z. Wu, M. Ismail, E. Serpedin, and K. A. Qaraqe, "Q-learning based two-timescale power allocation for multi-homing hybrid RF/VLC networks," *IEEE Wireless Commun. Lett.*, vol. 9, no. 4, pp. 443–447, Apr. 2020.

[11] P. Luo, M. Zhang, Z. Ghassemlooy, H. Le Minh, H.-M. Tsai, X. Tang, and D. Han, "Experimental demonstration of a 1024-QAM optical camera communication system," *IEEE Photon. Technol. Lett.*, vol. 28, no. 2, pp. 139–142, Oct. 2015.

[12] R. V. Terres, "Multi-user MISO for visible light communication," Ph.D. dissertation, University of Virginia, Sep. 2015.

[13] A. Islam, L. Musavian, and N. Thomos, "Performance analysis of vehicular optical camera communications: Roadmap to uRLLC," in *Proc. IEEE GlobeCom,19*, Hawaii, USA, Dec. 2019, pp. 1–6.

[14] A. Ashok, M. Gruteser, N. Mandayam, J. Silva, M. Varga, and K. Dana, "Challenge: Mobile optical networks through visual MIMO," in *Proc. MobiCom'10*, New York, NY, USA, Sep. 2010, pp. 105–112.

[15] B. Horn, B. Klaus, and P. Horn, *Robot vision*. MIT press, 1986.

[16] P. Deng, "Real-time software-defined adaptive MIMO visible light communications," *Visible Light Communications*, pp. 637–640, Jul. 2017.

[17] A. Ashok, S. Jain, M. Gruteser, N. Mandayam, W. Yuan, and K. Dana, "Capacity of screen–camera communications under perspective distortions," *Pervasive Mob. Comput.*, vol. 16, pp. 239–250, Jan. 2015.

[18] D. A. Plaisted, "Some polynomial and integer divisibility problems are NP-HARD," in *Proc. SFCS'76*, TX, USA, Oct. 1976, pp. 264–267.

[19] T. Zinchenko, "Reliability assessment of vehicle-to-vehicle communication," doctoralthesis, Technische Hochschule Wildau, 2014.

[20] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.

[21] N. Mastronarde and M. van der Schaar, "Joint physical-layer and system-level power management for delay-sensitive wireless communications," *IEEE Trans. Mob. Comput.*, vol. 12, no. 4, pp. 694–709, Feb. 2012.

[22] N. Salodkar, A. Bhorkar, A. Karandikar, and V. S. Borkar, "An on-line learning algorithm for energy efficient delay constrained scheduling over a fading channel," *IEEE J. Sel. Areas Commun.*, vol. 26, no. 4, pp. 732–742, Apr. 2008.

[23] D. Krajzewicz, J. Erdmann, M. Behrisch, and L. Bieker, "Recent development and applications of SUMO-simulation of urban mobility," *Int. J. Advances Syst. Measurements*, vol. 5, no. 3&4, Dec. 2012.

[24] M. Abadi *et al.*, "Tensorflow: A system for large-scale machine learning," in *Proc. OSDI '16*, Savannah, GA, Nov. 2016, pp. 265–283.