

Collaborative Data Scheduling for Vehicular Edge Computing via Deep Reinforcement Learning

Quyuan Luo, Changle Li, *Senior Member, IEEE*, Tom H. Luan, *Senior Member, IEEE*, and Weisong Shi *Fellow, IEEE*

Abstract—With the development of autonomous driving, the surging demand for data communications as well as computation offloading from connected and automated vehicles can be expected in the foreseeable future. With the limited capacity of both communication and computing, how to efficiently schedule the usage of resources in the network towards best utilization represents a fundamental research issue. In this paper, we address the issue by jointly considering the communication and computation resources for data scheduling. Specifically, we investigate on the Vehicular Edge Computing (VEC) in which edge computing enabled roadside unit (RSU) are deployed along the road to provide data bandwidth and computation offloading to vehicles. In addition, vehicles can collaborate among each other with data relays and collaborative computing via vehicle-to-vehicle (V2V) communications. A unified framework with communication, computation, caching, and collaborative computing is then formulated, and a collaborative data scheduling scheme to minimize the system-wide data processing cost with ensured delay constraints of applications is developed. To derive the optimal strategy for data scheduling, we further model the data scheduling as a deep reinforcement learning problem which is solved by an enhanced deep Q-network (DQN) algorithm with a separate target Q-network. Using extensive simulations, we validate the effectiveness of the proposal.

Index Terms—Vehicular edge computing; data scheduling; deep reinforcement learning; deep Q-network.

I. INTRODUCTION

IN RECENT years, there are increasing interest in connected and automated vehicles, which integrate information and communication technologies and play a crucial role toward a safer and more intelligent transportation system [1]. However, the increasing number of connected and automated vehicles and their resource hungry applications pose great challenges to the limited capability of vehicles in terms of data computing capacity for providing real-time and reliable vehicular services [2]. The paradigm of Vehicular Edge Computing (VEC) has been came up accordingly to strengthen the service capability through moving computation nodes to proximity of vehicles [3].

This work was supported by the National Natural Science Foundation of China under Grant No. U1801266, Key Research and Development Program of Shaanxi (Contract No. 2018ZDXM-GY-038, 2018ZDCXL-GY-04-02), the Youth Innovation Team of Shaanxi Universities, the Science and Technology Projects of Xian, China (201809170CX11JC12), and China Scholarship Council. (Corresponding author: Changle Li)

Q. Luo and C. Li are with the State Key Laboratory of Integrated Services Networks, Xidian University, Xi'an 710071, China (e-mail: qyluo@stu.xidian.edu.cn; clli@mail.xidian.edu.cn).

T. H. Luan is with the School of Cyber Engineering, Xidian University, Xi'an 710071, China (e-mail: tom.luan@xidian.edu.cn).

W. Shi is with the Department of Computer Science, Wayne State University, Detroit, MI 48202, USA (e-mail: weisong@wayne.edu).

In a VEC network, edge computational nodes (ECNs) can be deployed in cell towers, Road-Side Units (RSUs), and within connected and automated vehicles. Although offloading data to ECNs can significantly improve the delay performance, it also increases the burden of radio spectrum resource when data is transmitted through radio access networks [4], [5]. The sharp increase in the demand for communication during data offloading imposes great challenges for the existing communication resources of the VEC network [6]. Accordingly, it is vital to jointly consider the communication and computation resource for data scheduling to alleviate the situation.

Some existing works have focused on the joint allocation of communication and computation resources [7]–[10]. In these work, the data scheduling is mostly formulated as a resource scheduling problem through either minimizing the total delay or cost, or maximizing the system utility. Marvelous solutions are proposed to solve these optimization problems. However, the dynamic and changing VEC environment make the resource allocation an non-convex optimization problem with complicated objection function and constraints, which are much complicated and difficult to solve [11]. Good at abstracting many factors that affect the VEC environment to a mapping problem and learning optimal resource allocation strategy from the environment, deep reinforcement learning (DRL) is becoming a research hotspot to solve resource allocation problem [12], [13]. Tan *et al.* [14] propose a DRL-based joint communication, caching and computing allocation scheme in vehicle networks. He *et al.* [15] propose a dueling-DQN (DDQN)-based approach to solve their formulated joint optimization problem of networking, caching and computing resources for connected vehicles. More recently, Zhang *et al.* [16] propose a DQN empowered task offloading for MEC in Urban informatics.

However, most of the existing works consider offloading data to RSU or cell towers, ignoring the idle computing capacity of collaborative vehicles. Since vehicles themselves can be regarded as ECNs, they can assist data processing for other vehicles if they have some idle computing resources, which will greatly reduce the data processing cost and burden of RSUs or cell towers. In light of the existing works, in this paper, we make a further step in designing, analyzing and optimizing data scheduling through jointly consider communication, computation, caching, and collaborative computing a unified framework. Inspired by the deep reinforcement learning (DRL), we further propose a collaborative data scheduling scheme for VEC network based on deep Q-network (DQN), to minimize data processing cost while ensuring delay constraints. Specifically, the contributions of this paper can be

summarized as follows.

- 1) *Model*: We model the data scheduling model in a unified VEC network, where data can be processed locally, offloaded to RSUs, migrated to collaborative vehicles, or kept in caching queues. Considering the remaining lifetime and caching state of data, we establish a multi-queue model for data caching on both vehicle side and RSU side.
- 2) *Algorithm Design*: To derive the optimal data scheduling strategy, we first formulate an optimization problem to minimize the system-wide data processing cost while ensuring delay constraints. Then the data scheduling is modeled as a deep reinforcement learning problem to reflect the interaction between data scheduling and cost, by jointly considering communication resource, caching states of data, computing resource, delay requirements of data, and the mobility of vehicles.
- 3) *Validation*: Based on the real-world vehicular trace, the proposed scheme is evaluated by extensive simulations. Simulation results validate the performance of our proposal and show that our proposal efficiently reduces data processing cost and helps data be processed under delay constraints.

The remainder of this paper is organized as follows. The related work is presented in Section II. In Section III, we depict the system model and establish a multi-queue model for data caching on dual-sides. The data scheduling problem is formulated in Section IV. The DQN-based scheduling scheme is introduced in Section V. In Section VI, extensive simulation results are discussed. The conclusion is drawn in Section VII.

II. RELATED WORK

In this section, we survey the existing literature on joint resource allocation of communication and computation for data scheduling in vehicular edge computing system.

Du *et al.* [19] propose an online DDORV algorithm, which utilizes Lyapunov optimization theory to minimize the averaged cost of MEC enabled roadside units (MRSU) and vehicular terminal. Through derivation and comparing the values of local processing cost and task offloading cost, the optimization problem on vehicular terminal side is solved. For the optimization issue on MEC server side, the Lagrangian dual decomposition and continuous relaxation method is adopted. Zhang *et al.* [17] propose a cloud-based MEC offloading framework in vehicular networks, where both the heterogeneous requirements of the mobility of the vehicles and the computation tasks are considered. Based on the analysis of the characteristics of various offloading strategies, the authors further propose a predictive-mode transmission scheme for task-file uploading. Moreover, in [8], Zhang *et al.* adopt a Stackelberg game theory approach to design an optimal multilevel offloading scheme, which maximizes the utilities of both the vehicles and the computing servers.

Most of the existing MEC based resource allocation and optimization problems are mixed-integer non-linear programming (MINLP) problems and they are also NP-hard problems, which are not computable in polynomial time with existing

general time with existing general solvers [18]. Generally speaking, the complex optimization problem can be decomposed into subproblems, and by solving the subproblems respectively, the near-optimal solution is derived [7]–[9], [19], [20].

In order to efficiently solve the resource allocation optimization problem, some other methods such as game theory, heuristic intelligent algorithm, deep reinforcement learning are also becoming the research focuses. Messous *et al.* [21] consider the problem of computation offloading while achieving a tradeoff between execution time and energy consumption in an unmanned aerial vehicle (UAV) network, where the combination of energy overhead and delay is minimized by the designed game theory model. Dinh *et al.* [22] formulate a distributed computation offloading problem among the mobile users as an exact potential game and propose a distributed offloading scheme based on Q -learning and better-response with inertia and prove the Nash equilibrium convergence.

As to deep reinforcement learning based approach, Zhang *et al.* [2] adopt a deep Q -learning approach for designing an optimal data transmission scheduling scheme in cognitive vehicular networks to minimize transmission costs while also fully utilizing various communication modes and resources. In [16], they also design an optimal offloading scheme with joint MEC server selection and transmission mode determination in a deep Q -learning approach to maximize task offloading utility. Tan *et al.* [14] develop a framework of joint optimal resource allocation of communication, caching and computing for vehicular networks and propose a deep reinforcement learning approach to solve this resource allocation problem. He *et al.* [15] propose an integrated framework that can enable dynamic orchestration of networking, caching and computing resources for connected vehicles and formulate the resource allocation strategy as a joint optimization problem. To derive the optimal strategy, they proposed a dueling-DQN (DDQN)-based approach to solve the problem.

All those works above are marvelous solutions. Most existing relevant works focus on how to allocate resource of networking, caching and computing for better data scheduling and processing. The data is generally offloaded to RSUs or cell towers to compute. Few works focus on the collaborative computing on vehicle side. This work takes the research a step further by fully utilizing the idle computing resources of collaborative vehicles, formulating a unified framework with communication, computation, caching, and collaborative computing, developing a collaborative data scheduling scheme to minimize the system-wide data processing cost with ensured delay constraints of applications.

III. SYSTEM MODEL

This section presents the system model including the unified framework with communication, computation, caching, and collaborative computing; and the multi-queue model for data caching on dual-sides. For convenience, the main notations used are summarized in Table I.

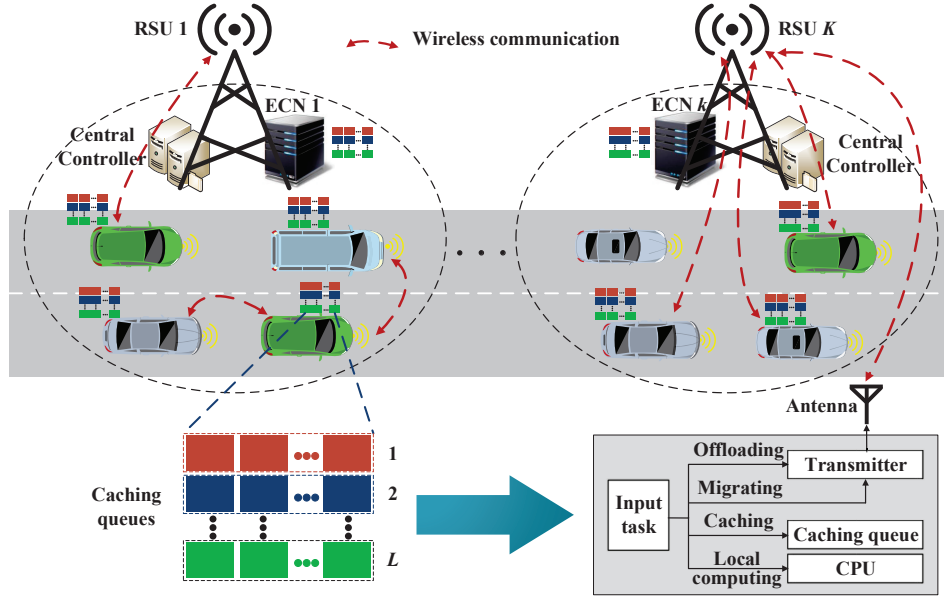


Fig. 1: Unified framework with communication, computation, caching, and collaborative computing in VEC network.

TABLE I: Major Notations

Notation	Explanation
K	Number of road segments
C_{V2I}	Number of licensed channels for V2I communication
C_{V2V}	Number of licensed channels for V2V communication
B	Bandwidth of each licensed channel
R_k	Coverage radius of RSU k
R^V	Coverage radius of vehicles
N_k	Number of vehicles within RSU k
Δt	Duration of a time-slot
M	Number of data types
D_i	Amount of type- i data
c	Processing density of data
f_n^{local}	Processing capability of vehicle n
f_k^{ECN}	Processing capability of RSU k
κ_1, κ_2	Effective switched capacitance related to the chip architecture in vehicles and RSUs
d	Distance between transmitter and receiver
ϑ	Path loss exponent
h	Channel fading coefficient
ω_0	White Gaussian noise power
P_n^{tr}	Transmission power of vehicle n
$\alpha_{n,k}^t, \beta_{n,k}^t$	Local computing and data offloading indicators for vehicle n running on road segment k at time-slot t
$\gamma_{n,k}^t, \delta_{n,k}^t$	Data migrating and receiving indicators for vehicle n running on road segment k at time-slot t
μ_n^t	Data processing indicator for RSU n at time-slot t
ξ	Penalty coefficient
$q_{n,l}^t, g_{k,l}^t$	Length of data cached in queue l of vehicle n and RSU k at time-slot t
c^I, c^V	Cost for using licensed V2I and V2V channels
c^{ECN}	Cost for RSU processing data
ϱ	Cost for energy consumption

A. Unified Framework with Communication, Computation, Caching, and Collaborative Computing in VEC Network

Fig. 1 shows the unified framework with communication, computation, caching, and collaborative computing in VEC network. The road is divided into K segments, denoted by $\mathbb{K} = \{1, 2, \dots, K\}$, and each covered by a roadside unit (RSU) with an edge computing node (ECN). The coverage radius

of these RSUs are denoted by $\{R_1, R_2, \dots, R_K\}$, respectively. Vehicles can establish communication link with both RSU and vehicles through the assigned orthogonal licensed channels, each with the bandwidth of B . The licensed bandwidth is classified into two categories, one is for vehicle-to-infrastructure (V2I) communication, and the other for vehicle-to-vehicle (V2V) communication. The number of licensed channel for V2I and V2V communications are denoted by C_{V2I} and C_{V2V} , respectively. And RSUs also provide power computing capacity due to the deployed ECN.

In a VEC scenario, various data would be generated from on-board applications both for safety (e.g., high-definition camera and LiDAR) and entertainment (e.g., augmented reality and face recognition) purpose [23], [24]. These applications generally utilize deep learning method to process data, which need powerful computational capacity [25]. In order to better describe the data generating, transmitting and computing processes, we divide time into time-slots, each with length of Δt . Since Δt is short, we consider the system be quasi-static so that the wireless channels and the topology of the system keep unchanged at each time-slot and vary at different time-slots [19]. And data would be generated at the beginning of each time-slot.

Due to the limited computing capacity of on-board device, severe delay would be caused for processing computation-intensive data. Thanks to the licensed channel resources, some latency-sensitive and computation-intensive data of vehicles could be offloaded to RSUs or migrated to collaborative vehicles that have idle computation resources. The control center in each road segment schedules vehicular communications through a dedicated control channel. Besides, the control center plays a role in gathering the states of vehicles through pilot signals. The size of pilot signal can be set to arbitrarily small. Thus, the extra overhead can be very small [26]. Since processing result is usually very tiny, we neglect the output return process and just focus on transmitting data to RSUs or

collaborative vehicles [19]. We categorize data into M types and use two items $\{D_i, T_i\}$ to describe an arbitrary type of data i ($i \in \mathbb{M} = \{1, 2, \dots, M\}$), where D_i stands for the amount of data, and T_i stands for the data delay constraint. In addition, we use c to represent the processing density (in CPU cycles/bit) of the data. It is noteworthy that the unit of delay constraint is time-slot, which means once the data is generated, it needs to be processed within time duration of $\Delta t \times T_i$. At each time-slot, type- i data is generated with the probability of ϖ_i , which meets $\sum_{i=1}^M \varpi_i \leq 1$. The generated data can be processed locally, or offloaded to a RSU or migrated to a collaborative vehicle to be processed. For some delay-tolerant data, it may give way to the delay-sensitive data in the vehicular caching for priority processing. In order to better illustrate it, we model the data caching on both vehicle side and RSU side as a multi-queue system in the following.

B. Multi-Queue Model for Data Caching on Dual-Sides

For the data caching on both vehicle side and RSU side, we use L to denote the number of caching queues indexed as $\{1, 2, \dots, L\}$, respectively. According to the delay constraints of generated data, we have $L = \max\{T_i, i \in \mathbb{M}\}$. In each queue, the remaining lifetime of data under its delay constraint is the same. Since the caching queue's input and output of vehicles and RSUs are different, we analyze them respectively in the following.

a) Vehicle Side: On the vehicle side as shown in Fig. 2(a), we divide the queue into two categories according to the queue index l , i.e., $1 \leq l < L$ and $l = L$. Assume the index of time-slot now is t , for the case when $1 \leq l < L$, the input data in caching queue l in a given vehicle V1 comes from three sources:

- 1) The data in $l + 1$ of V1 and has not been processed at time-slot $t - 1$. As time-slot increases from $t - 1$ to t , the remaining lifetime for data processing decreases by 1.
- 2) The data generated by vehicle V1 itself at time-slot t with data delay constraint l .
- 3) Another vehicle V2 transmitted data with queue index $l + 1$ to V1 at time-slot $t - 1$.

The data in queue l of vehicle V1 when $1 \leq l < L$ has four different outputs:

- 1) Be offloaded to RSU k and then cached in queue $l - 1$ of RSU k at next time-slot $t + 1$.
- 2) Be migrated to another vehicle V3 and then cached in $l - 1$ of V3 at next time-slot $t + 1$.
- 3) Be computed locally through the computation resource of vehicle V1.
- 4) If none of the three actions above is adopted, data in queue l will be moved to queue $l - 1$ of vehicle V1 if $l \neq 1$, otherwise deleted if $l = 1$, at next time-slot $t + 1$.

For the case when $l = L$, the input data in caching queue l only comes from the newly generated data with delay constraint L . And the output of data in queue l when $l = L$ has four transitions, the same as when $1 \leq l < L$.

b) RSU Side: On the RSU side, as shown in Fig. 2(b), we also detail the input and output of data in each queue l . We

divide the queue into two categories according to the queue index, i.e., $1 \leq l < L - 1$ and $l = L - 1$. It is noteworthy that the maximal queue index l of RSUs is $L - 1$. This is because data can not be generated by RSUs themselves and it will take at least one time-slot for data to be offloaded from vehicles to RSUs. Just like the analysis on the vehicle side, we assume the index of time-slot now is t , for the case when $1 \leq l < L - 1$, the input data in caching queue l in a given RSU k comes from two sources¹:

- 1) The data in $l + 1$ of RSU k and has not been processed at time-slot $t - 1$.
- 2) Vehicles transmitted data with queue index $l + 1$ to RSU k at time-slot $t - 1$.

The data in queue l of RSU k when $1 \leq l < L - 1$ has two different outputs:

- 1) Be computed through the computation resource of ECN-enable RSU.
- 2) Be moved to queue $l - 1$ of RSU k if $l \neq 1$, otherwise deleted if $l = 1$, at next time-slot $t + 1$.

For the case when $l = L - 1$, the input data in caching queue l only consist of the data with remaining lifetime $l + 1$ (i.e., data with remaining lifetime L) offloaded from vehicles through V2I communication at time-slot $t - 1$. And the output of data in queue l when $l = L - 1$ has two different outputs, the same as when $1 \leq l < L - 1$.

IV. DATA SCHEDULING: ANALYSIS, PROBLEM FORMULATION, AND MDP

In this section, we first investigate the performance of the formulated unified VEC framework with various transmission and computation modes. Then we formulate an optimal data scheduling problem that take into account both cost and delay constraints.

A. Analysis of Data Scheduling: Transmission and Computation

The data cached in the caching queues of both vehicles and RSUs can be scheduled in different ways. We will detail them in the following.

1) Data processed locally: We denote the processing capability (i.e., the amount of CPU frequency in cycles/s) at vehicle n assigned for local computing as f_n^{local} , then the power consumption for vehicle n to process data locally is expressed as

$$p_n^{\text{local}} = \kappa_1 (f_n^{\text{local}})^3, \quad (1)$$

where κ_1 stands for the effective switched capacitance related to the chip architecture in vehicle [27]. Accordingly, the energy consumption of vehicle n for local processing during one time-slot is expressed as

$$E_n^{\text{local}} = p_n^{\text{local}} \Delta t = \kappa_1 (f_n^{\text{local}})^3 \Delta t. \quad (2)$$

And the amount of data be processed locally by vehicle n during one time-slot is expressed as

$$D_n^{\text{local}} = \frac{f_n^{\text{local}} \Delta t}{c}. \quad (3)$$

¹In this paper, we do not consider the migration of data between RSUs, which will be considered as our future work.

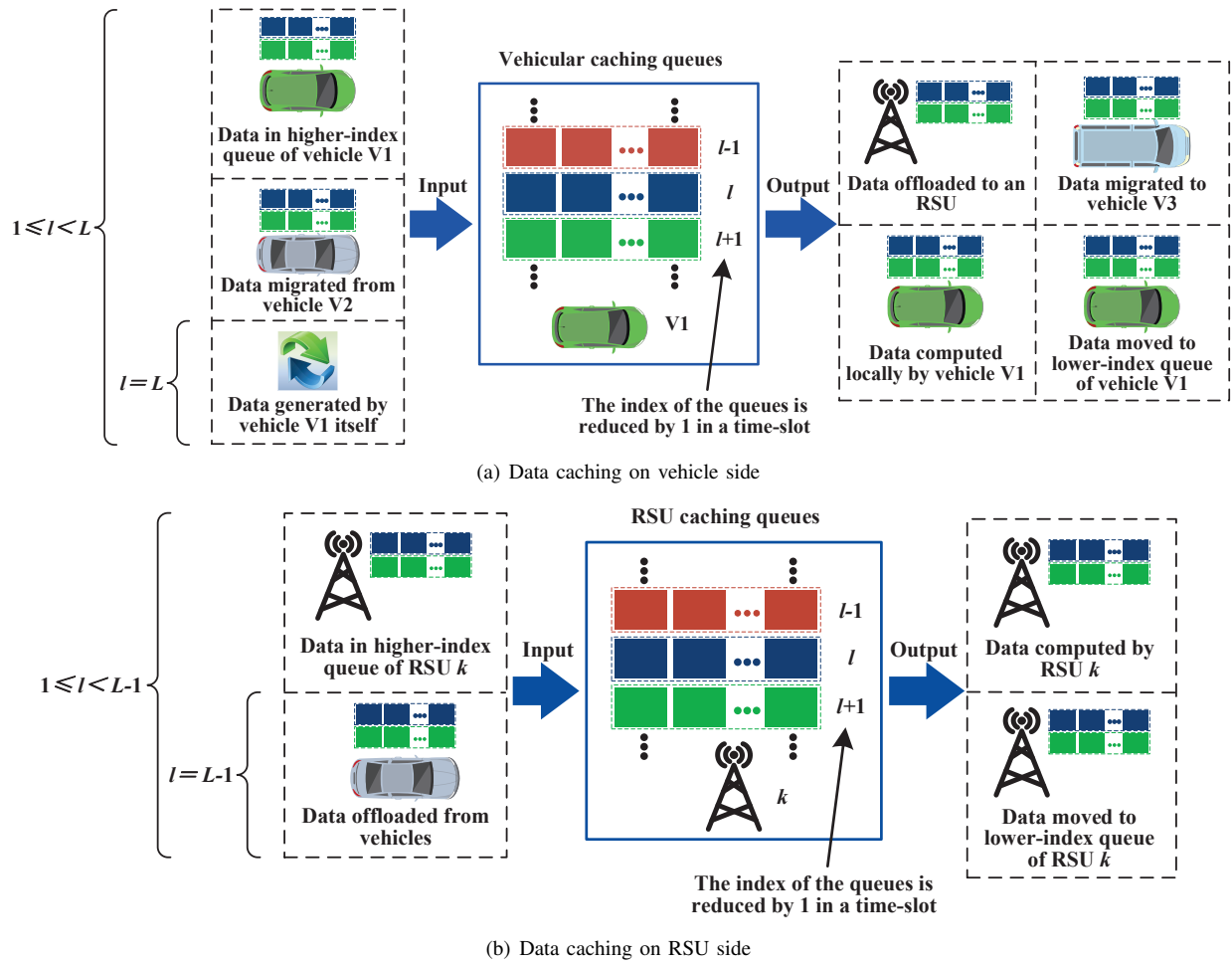


Fig. 2: Multi-queue model for caching data on dual-sides: a) multi-queue model for data caching on vehicle side. According to the queue index, there are two cases, i.e., when $1 \leq l < L$ and when $l = L$. For $1 \leq l < L$, the input data comes from three sources. For $l = L$, the input data only comes from one source. b) multi-queue model for data caching on RSU side.

2) *Data offloaded to RSU:* In VEC network, the data can be offloaded to RSU through V2I communication scheduled by the control center. To model the data offloading through V2I communication, we denote the path loss as $d^{-\vartheta}$, where d and ϑ denote the distance from transmitter to receiver and the path loss exponent, respectively. Moreover, the channel fading coefficient is denoted by h , which is modeled as a circularly symmetric complex Gaussian random variable [28]. When data is offloaded from vehicle n to RSU k on a licensed V2I channel, the transmission rate is given by

$$r_{n,k}^{t,V2I} = B \log_2 \left(1 + \frac{P_n^{\text{tr}} |h|^2}{\omega_0 (d_{n,k}^t)^{\vartheta}} \right), \quad (4)$$

where P_n^{tr} is the transmission power of vehicle n , ω_0 denotes the white Gaussian noise power, $d_{n,k}^t$ denotes the distance from vehicle n to RSU k at time-slot t . The energy consumption of vehicle n for transmitting data during one time-slot is expressed as

$$E_n^{\text{tr}} = P_n^{\text{tr}} \Delta t. \quad (5)$$

3) *Data processed by RSU:* The data in the caching queues of RSUs is processed by deployed ECNs. We denote the

processing capability at RSU k as f_k^{ECN} , then the power consumption for RSU k to process data is expressed as

$$p_k^{\text{ECN}} = \kappa_2 (f_k^{\text{ECN}})^3, \quad (6)$$

where κ_2 stands for the effective switched capacitance related to the chip architecture in RSU [27]. Accordingly, the energy consumption of RSU k for data processing during one time-slot is expressed as

$$E_k^{\text{ECN}} = p_k^{\text{ECN}} \Delta t = \kappa_2 (f_k^{\text{ECN}})^3 \Delta t. \quad (7)$$

And the amount of data processed by RSU k during one time-slot is expressed as

$$D_k^{\text{ECN}} = \frac{f_k^{\text{ECN}} \Delta t}{c}. \quad (8)$$

4) *Data migrated to collaborative vehicle:* When some vehicles have no data to process or their data in caching queues is delay-tolerant and has longer remaining lifetime, they can aid data processing as collaborative vehicles. In this way, data can be also migrated to collaborative vehicles through V2V communication. We denote the communication radius of vehicles as R^V . For arbitrary vehicle n and collaborative

vehicle n' , n can transmit data to n' while n does not take any transmission action at the same time-slot. The communication model is similar to V2I communication model and formula (4). When data is migrated from n to n' on a licensed V2V channel at time-slot t , the transmission rate is given by

$$r_{n,n'}^{t,V2V} = B \log_2 \left(1 + \frac{P_n^{\text{tr}} |h|^2}{\omega_0 (d_{n,n'}^t)^{\bar{\theta}}} \right), \quad (9)$$

where $d_{n,n'}^t$ is the distance between vehicle n and n' . It is noteworthy that only when $d_{n,n'}^t < R^V$ does formula (9) hold.

B. Problem Formulation

At a give time-slot, the data of each vehicle can be: a) processed locally; b) offloaded to RSUs; c) migrated to collaborative vehicles; and d) kept in its caching queues. Let $\alpha_{n,k}^t = 1$ indicate that the data of vehicle n running on road segment k at time-slot t is processed locally. Similarly, we use $\beta_{n,k}^t = 1$ to indicate that the data of vehicle n is offloaded to RSU through V2I communication, $\gamma_{n,k}^t = 1$ to indicate that the data of vehicle n is migrated to a collaborative vehicle through V2V communication. The case $\alpha_{n,k}^t = \beta_{n,k}^t = \gamma_{n,k}^t = 0$ indicates that vehicle n keeps the data in its caching queues. In addition, vehicle n can receive data migrated from other vehicles, let $\delta_{n,k}^t$ indicate that vehicle n on road segment k is on receiving mode at time-slot t . On the RSU side, each RSU can process data or keep the data in its caching queues. We use $\mu_k^t = 1$ to indicate that RSU k processes data at time-slot t , $\mu_k^t = 0$ to indicate that RSU k keeps the data in its caching queues.

To promote the data processing, we introduce a penalty mechanism, which means a penalty will be resulted in if data is not processed before its deadline. We use ξ to denote the penalty coefficient indicating the penalty amount of one unit data. Moreover, costs for communication and computation are also produced during data scheduling. Those costs includes: a) the cost for using licensed channels of V2I and V2V communications; b) the cost for RSU computing data; and c) the energy consumption cost during data computing and transmitting. To be able to process data under delay constraints, the data with smaller queue index should be processed firstly. We use \mathcal{V}_k^t to denote vehicle set in road segment k at time-slot t . The amount of data not meeting the delay constraints at time-slot t is then expressed as

$$D_{\text{loss}}^t = \sum_{k \in \mathbb{K}} \left(\sum_{n \in \mathcal{V}_k^t} \alpha_{n,k}^t \max \left\{ 0, q_{n,1}^t - D_n^{\text{local}} \right\} + \mu_k^t \max \left\{ 0, g_{k,1}^t - D_k^{\text{ECN}} \right\} + (1 - \alpha_{n,k}^t) q_{n,1}^t + (1 - \mu_k^t) g_{k,1}^t \right), \quad (10)$$

where $q_{n,1}^t$ and $g_{k,1}^t$ denote the amount of data cached in queue 1 of vehicle n and RSU k at time-slot t , respectively. D_n^{local} and D_k^{ECN} are the data processing capabilities of vehicles and RSUs at one time-slot, as formulated by (3) and (8), respectively.

To reduce cost while ensuring delay constraints through making full use of communication and computation resources, we formulate the data scheduling problem as

$$\begin{aligned} \min_{\{\alpha, \beta, \gamma, \delta, \mu\}} \text{Loss} = & \sum_{t=1}^{\infty} \left\{ \xi D_{\text{loss}}^t + \sum_{k \in \mathbb{K}} \left(\varrho E_n^{\text{local}} \sum_{n \in \mathcal{V}_k^t} \alpha_{n,k}^t \right. \right. \\ & + (c^l + \varrho E_n^{\text{tr}}) \sum_{n \in \mathcal{V}_k^t} \beta_{n,k}^t (1 - \delta_{n,k}^t) \\ & + (c^V + \varrho E_n^{\text{tr}}) \sum_{n \in \mathcal{V}_k^t} \gamma_{n,k}^t (1 - \delta_{n,k}^t) \\ & \left. \left. + (c^{\text{ECN}} + \varrho E_k^{\text{ECN}}) \mu_k^t \right) \right\} \\ \text{s.t. } & \text{C1: } \alpha_{n,k}^t, \beta_{n,k}^t, \gamma_{n,k}^t, \delta_{n,k}^t, \mu_k^t \in \{0, 1\}, \\ & \text{C2: } \beta_{n,k}^t \gamma_{n,k}^t = 0, \\ & \text{C3: } \beta_{n,k}^t \delta_{n,k}^t = \gamma_{n,k}^t \delta_{n,k}^t = 0, \\ & \text{C4: } \alpha_{n,k}^t \beta_{n,k}^t = \alpha_{n,k}^t \gamma_{n,k}^t = 0, \end{aligned} \quad (11)$$

where ϱ is a weight coefficient indicating the energy consumption cost of one unit energy during data computing and transmitting [29], c^l and c^V denote the costs at a time-slot for using licensed channels for V2I and V2V communications, respectively. And c^{ECN} denotes the cost for RSU processing data at a time-slot. In formula (11), constraint C1 indicates that whether a vehicle or a RSU takes one action or not at time-slot t . Constraint C2 indicates that vehicles cannot transmit data through V2I and V2V communications simultaneously. Constraint C3 indicates that when a vehicle is on receiving mode, it cannot transmit data neither through V2I nor V2V communication. Constraint C4 indicates that at most one data processing action among local computing, offloading through V2I, and migrating through V2V can be chosen during one time-slot.

C. Model Data Processing Scheduling as a MDP

In formula (11), *Loss* mainly depends on the states and data scheduling actions of both vehicles and RSUs. And states of next time-slot only depends on the current state and the data scheduling actions. Accordingly, the data scheduling problem can be formulated as a Markov decision process (MDP) to analysis the state transitions of caching queues.

We denote the state of the MDP at time-slot t as $S^t \triangleq \{S_1^t, S_2^t, \dots, S_K^t, \Phi^t\}$, where Φ^t is the position state of vehicles and RSUs, S_k^t ($1 \leq k \leq K$) is the caching state of vehicles and RSU in road segment k . Since the positions of RSUs are fixed and the next positions of vehicles can be obtained based on the current positions and running speeds, the next state Φ^{t+1} can be easily obtained. Therefore, we mainly focus on the caching state transitions. In arbitrary road segment k , we define $S_k^t \triangleq \{Q_{1,k}^t, Q_{2,k}^t, \dots, Q_{|\mathcal{V}_k^t|,k}^t, G_k^t\}$, where $|\mathcal{V}_k^t|$ is the number of vehicles in road segment k at time-slot t , $Q_{n,k}^t$ ($n \in \mathcal{V}_k^t$) is the caching state of vehicle n and is expressed as $\{q_{n,1}^t, q_{n,2}^t, \dots, q_{n,L}^t\}$, G_k^t is the caching state of RSU k and is expressed as $\{g_{k,1}^t, g_{k,2}^t, \dots, g_{k,L}^t\}$. $q_{n,l}^t$ and $g_{k,l}^t$ ($1 \leq l \leq L$) denote the amount of data cached in queue l of

vehicle n and RSU k at time-slot t , respectively. We define the action taken by RSUs and vehicles as $A^t \triangleq \{A_1^t, A_2^t, \dots, A_K^t\}$. And for the RSU and vehicles in road segment k , we define $A_k^t \triangleq \{a_{1,k}^t, a_{2,k}^t, \dots, a_{|\mathcal{V}_k^t|,k}^t, \tilde{a}_k^t\}$, where $a_{n,k}^t$ ($n \in \mathcal{V}_k^t$) and \tilde{a}_k^t are the actions of vehicle n and RSU k , respectively. $a_{n,k}^t$ is consisted of the possible data scheduling actions, expressed as $a_{n,k}^t = \{\alpha_{n,k}^t, \beta_{n,k}^t, \gamma_{n,k}^t, \delta_{n,k}^t\}$. \tilde{a}_k^t can be expressed as $\{\mu_k^t\}$. To describe the caching state at the next time-slot $t+1$, we should calculate the amount of data transmission during the current time-slot t . The amount of data that is transmitted from vehicle n to RSU k can be expressed as

$$D_{n,k}^{t,\text{off}} = \begin{cases} \frac{\beta_{n,k}^t r_{n,k}^{t,V2I} \Delta t C_{V2I}}{|\mathcal{V}_k^{t,\text{off}}|}, & \mathcal{V}_k^{t,\text{off}} \neq \emptyset \\ 0, & \mathcal{V}_k^{t,\text{off}} = \emptyset \end{cases} \quad (12)$$

where $\mathcal{V}_k^{t,\text{off}}$ denotes the set of vehicles in road segment k choosing to offload data to RSU at time-slot t . In addition, we use $\mathcal{V}_k^{t,\text{mig}}$ and $\mathcal{V}_k^{t,\text{rec}}$ to denote the sets of vehicles in road segment k choosing to transmit data to collaborative vehicles and choosing to receive migration data at time-slot t , respectively. We define an indicator $I_{n,n'}^{t,k}$ to denote the connection between vehicle n ($n \in \mathcal{V}_k^{t,\text{mig}}$) and n' ($n' \in \mathcal{V}_k^{t,\text{rec}}$), which equals 1 when the connection is established and 0 otherwise, $I_{n,n'}^{t,k}$ is defined as

$$I_{n,n'}^{t,k} = \begin{cases} 1, & \gamma_{n,k}^t = 1, \delta_{n',k}^t = 1, d_{n,n'} < R^V \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

We assume that a vehicle n can only transmit data to at most one collaborative vehicle at one time-slot, hence $I_{n,n'}^{t,k}$ meets the condition that $\sum_{n' \in \mathcal{V}_k^{t,\text{rec}}} I_{n,n'}^{t,k} \leq 1$. And the amount of data that is transmitted from vehicle n to n' can be expressed as

$$D_{n,n',k}^{t,\text{mig}} = \begin{cases} I_{n,n'}^{t,k} \frac{r_{n,n'}^{t,V2V} \Delta t C_{V2V}}{|\mathcal{V}_k^{t,\text{mig}}|}, & \mathcal{V}_k^{t,\text{mig}} \neq \emptyset \\ 0, & \mathcal{V}_k^{t,\text{mig}} = \emptyset \end{cases} \quad (14)$$

We use $l_{\min}^{t,n}$ to denote the smallest index of the queue with nonempty queueing data of arbitrary vehicle n at time-slot t . Then the amount of data migrated from queue $l+1$ of vehicle \bar{n} to queue l of vehicle n and RSU k at time-slot t are expressed as

$$D_{\bar{n},n,k,l}^{t,V2V} = \begin{cases} \mathbf{1}(l_{\min}^{t,\bar{n}} == l+1) \times \\ I_{\bar{n},n}^{t,k} r_{\bar{n},n}^{t,V2V} \Delta t C_{V2V} / |\mathcal{V}_k^{t,\text{mig}}|, & \mathcal{V}_k^{t,\text{mig}} \neq \emptyset \\ 0, & \mathcal{V}_k^{t,\text{mig}} = \emptyset \end{cases} \quad (15)$$

and

$$D_{\bar{n},k,l}^{t,V2I} = \begin{cases} \mathbf{1}(l_{\min}^{t,\bar{n}} == l+1) \times \\ \beta_{\bar{n},k}^t r_{\bar{n},k}^{t,V2I} \Delta t C_{V2I} / |\mathcal{V}_k^{t,\text{off}}|, & \mathcal{V}_k^{t,\text{off}} \neq \emptyset \\ 0, & \mathcal{V}_k^{t,\text{off}} = \emptyset \end{cases} \quad (16)$$

respectively, where $\mathbf{1}(\tau)$ is an indicator function which equals 1 if τ is true and 0 otherwise. Accordingly, given caching state $Q_{n,k}^t$ of vehicle n in road segment k , the state of queues that

form $Q_{n,k}^{t+1}$ is consisted of two parts, namely $q_{n,l_{\min}^{t,n}}^{t+1}$ and $q_{n,l}^{t+1}$ ($l \neq l_{\min}^{t,n}$), which are expressed as

$$q_{n,l_{\min}^{t,n}}^{t+1} = \begin{cases} \max \left\{ 0, q_{n,2}^t + \varpi_1 D_1 \right. \\ \quad \left. + \sum_{\bar{n} \in \mathcal{V}_k^{t,\text{mig}}} D_{\bar{n},n,k,1}^{t,V2V} \right. \\ \quad \left. - D_n^{\text{local}} \right\}, & l_{\min}^{t,n} = 1 \\ \max \left\{ 0, q_{n,l_{\min}^{t,n}+1}^t + \varpi_{l_{\min}^{t,n}} D_{l_{\min}^{t,n}}^{t,n} \right. \\ \quad \left. + \sum_{\bar{n} \in \mathcal{V}_k^{t,\text{mig}}} D_{\bar{n},n,k,l_{\min}^{t,n}}^{t,V2V} \right. \\ \quad \left. - \sum_{n' \in \mathcal{V}_k^{t,\text{rec}}} D_{n,n',k}^{t,\text{mig}} \right. \\ \quad \left. - D_{n,k}^{t,\text{off}} - D_n^{\text{local}} \right\}, & 1 < l_{\min}^{t,n} < L \\ \max \left\{ 0, \varpi_L D_L - D_{n,k}^{t,\text{off}} - D_n^{\text{local}} \right. \\ \quad \left. - \sum_{n' \in \mathcal{V}_k^{t,\text{rec}}} D_{n,n',k}^{t,\text{mig}} \right\}, & l_{\min}^{t,n} = L \end{cases} \quad (17)$$

and

$$q_{n,l}^{t+1} = \begin{cases} q_{n,l+1}^t + \varpi_l D_l \\ \quad + \sum_{\bar{n} \in \mathcal{V}_k^{t,\text{mig}}} D_{\bar{n},n,k,l}^{t,V2V}, & l \neq l_{\min}^{t,n}, 1 \leq l < L \\ \varpi_l D_l, & l \neq l_{\min}^{t,n}, l = L \end{cases} \quad (18)$$

respectively. Similarly, we use $\tilde{l}_{\min}^{t,k}$ to denote the smallest index of the queue with nonempty queueing data of arbitrary RSU k at time-slot t . Based on caching state G_k^t of RSU k , the state of queues that form G_k^{t+1} is consisted of two parts, namely $g_{k,\tilde{l}_{\min}^{t,k}}^{t+1}$ and $g_{k,l}^{t+1}$ ($l \neq \tilde{l}_{\min}^{t,k}$), which are expressed as

$$g_{k,\tilde{l}_{\min}^{t,k}}^{t+1} = \begin{cases} \max \left\{ 0, g_{k,\tilde{l}_{\min}^{t,k}+1}^t - D_k^{\text{ECN}} \right\} \\ \quad + \sum_{\bar{n} \in \mathcal{V}_k^{t,\text{off}}} D_{\bar{n},k,\tilde{l}_{\min}^{t,k}}^{t,V2I}, & 1 \leq \tilde{l}_{\min}^{t,k} < L-1 \\ \sum_{\bar{n} \in \mathcal{V}_k^{t,\text{off}}} D_{\bar{n},k,\tilde{l}_{\min}^{t,k}}^{t,V2I} - D_k^{\text{ECN}}, & \tilde{l}_{\min}^{t,k} = L-1 \end{cases} \quad (19)$$

and

$$g_{k,l}^{t+1} = \begin{cases} g_{k,l+1}^t + \sum_{\bar{n} \in \mathcal{V}_k^{t,\text{off}}} D_{\bar{n},k,l}^{t,V2I}, & l \neq \tilde{l}_{\min}^{t,k}, 1 \leq l < L-1 \\ \sum_{\bar{n} \in \mathcal{V}_k^{t,\text{off}}} D_{\bar{n},k,l}^{t,V2I}, & l \neq \tilde{l}_{\min}^{t,k}, l = L-1 \end{cases} \quad (20)$$

respectively.

In state S^t , the penalty and cost by taking action A^t can be expressed as

$$\begin{aligned} \text{Loss}^t = & \xi D_{\text{loss}}^t + \sum_{k \in \mathbb{K}} \left(\varrho E_n^{\text{local}} \sum_{n \in \mathcal{V}_k^t} \alpha_{n,k}^t \right. \\ & + (c^l + \varrho E_n^{\text{tr}}) \sum_{n \in \mathcal{V}_k^t} \beta_{n,k}^t (1 - \delta_{n,k}^t) \\ & + (c^V + \varrho E_n^{\text{tr}}) \sum_{n \in \mathcal{V}_k^t} \gamma_{n,k}^t (1 - \delta_{n,k}^t) \\ & \left. + (c^{\text{ECN}} + \varrho E_k^{\text{ECN}}) \mu_k^t \right). \end{aligned} \quad (21)$$

The purpose of the MDP is deriving an optimal data scheduling strategy that minimizes the cumulative value of $Loss^t$ over time-slots. The optimal strategy that indicates the data scheduling actions, is expressed as

$$\pi^* = \arg \min_{\pi} E \left(\sum_{t=1}^{\infty} \eta^t Loss^t \right), \quad (22)$$

where $0 < \eta < 1$ is a discount factor used to indicate the impact of future $Loss$ on current actions.

V. DQN-BASED OPTIMAL DATA SCHEDULING SCHEME

A. From Q -Learning to Deep Q -Network

In the formulated MDP problem (22), the large scale of state space and action space make it hard to find the optimal data scheduling strategy π^* [30]. Fortunately, the reinforcement learning technology is powerful in handling scheduling problem [13]. Reinforcement learning is a main branch of machine learning, where agents learn to take series of actions that maximize cumulative future reward with corresponding policies over states [15], [16]. Accordingly, our proposed MDP can be considered as a reinforcement learning problem, to minimize cumulative future loss. Under a given data scheduling strategy π , the expected long-term loss from taken action A^t at state S^t can be expressed as an action-value function (i.e., Q -function), which is shown as

$$\begin{aligned} Q_{\pi}(S^t, A^t) &= E \left[\sum_{i=0}^{\infty} \eta^i Loss^{t+i} \middle| (S^t, A^t) \right] \\ &= E \left[Loss^t + \eta^1 Loss^{t+1} + \dots \middle| (S^t, A^t) \right] \\ &= E_{S^{t+1}} \left[Loss^t + \eta Q_{\pi}(S^{t+1}, A^{t+1}) \middle| (S^t, A^t) \right]. \end{aligned} \quad (23)$$

When given action A^t in state S^t , the expected minimum loss is expressed as

$$\begin{aligned} Q^*(S^t, A^t) &= \\ E_{S^{t+1}} \left[Loss^t + \eta \min_{A^{t+1}} Q(S^{t+1}, A^{t+1}) \middle| (S^t, A^t) \right]. \end{aligned} \quad (24)$$

Based on formula (24), the minimum $Q^*(S^t, A^t)$ and optimal data scheduling actions can be derived by value and action iteration. The updated process of $Q(S^t, A^t)$, namely the Q -learning process, is expressed as

$$\begin{aligned} Q(S^t, A^t) &\leftarrow Q(S^t, A^t) + \varphi \left[Loss^t \right. \\ &\quad \left. + \eta \min_{A^{t+1}} Q(S^{t+1}, A^{t+1}) - Q(S^t, A^t) \right], \end{aligned} \quad (25)$$

where φ is the learning rate.

Since a Q -table is used in the Q -learning process to store learned state-action combinations and corresponding Q -values,

discrete state space is utilized in the Q -table. However, the states of the VEC network consist of the amount of data cached in the queues of vehicles and RSUs, whose value is continuous. Thus, Q -learning approach cannot be directly implemented in solving our proposed MDP problem. To compensate the limitation of Q -learning, we incorporate deep learning technology with Q -learning method, which forms the deep Q -network (DQN). Instead of Q -function, DQN uses a deep neural network as a nonlinear approximator that is able to capture the complex interaction among various states and actions. The inputs of the deep neural network are states, and the outputs are Q -values of actions. With the help of DQN, the Q -value in (23) can be estimated as $Q(S^t, A^t) \approx Q(S^t, A^t; \theta)$, where θ are the weights of the DQN. Accordingly, the optimal action for data scheduling in state S^t is the one with the minimum $Q(S^t, A^t; \theta)$, which is shown as

$$A^{t*} = \arg \min_{A^t} Q(S^t, A^t; \theta). \quad (26)$$

B. DQN Training

To guarantee the approximation ability of the estimated Q -value, $Q(S^t, A^t; \theta)$ should be trained towards the target value $Loss^t + \eta \min_{A^{t+1}} Q(S^{t+1}, A^{t+1})$, which is substituted with approximate target value as

$$y^t = Loss^t + \eta \min_{A^{t+1}} Q(S^{t+1}, A^{t+1}; \theta^{t-1}). \quad (27)$$

To minimize the difference between estimated and target values, we define a loss function as

$$Er(\theta^t) = E \left[\left(y^t - Q(S^t, A^t; \theta^t) \right)^2 \right], \quad (28)$$

where θ^t denotes the weights of the DQN at time-slot t . The parameters from the previous time-slot θ^{t-1} are held fixed when optimizing the loss function $Er(\theta^t)$ [31]. Differentiating the loss function with respect to the weights, we get the gradient as

$$\nabla_{\theta^t} Er(\theta^t) = E \left[2 \left(Q(S^t, A^t; \theta^t) - y^t \right) \nabla_{\theta^t} Q(S^t, A^t; \theta^t) \right]. \quad (29)$$

Based on gradient descent, θ^t is updated as

$$\theta^t \leftarrow \theta^t - \varepsilon \nabla_{\theta^t} Er(\theta^t), \quad (30)$$

where ε is a step size coefficient, which controls the updating step size in each iteration.

In order to improve learning efficiency while removing the correlations in the subsequent training samples at the same time, experience replay technique is utilized in the learning process. The learned experience $e^t = (S^t, A^t, Loss^t, S^{t+1})$ at each time-slot is stored in a data-set \mathcal{D} in a replay memory [32]. Then we randomly draw a batch of stored experience as samples to train parameters of DQN. Before performing experience replay, an action is selected and executed according to an ϵ -greedy policy, which avoids local optimum while balancing exploration and exploitation during training. That is, a random action is chosen with probability ϵ to explore

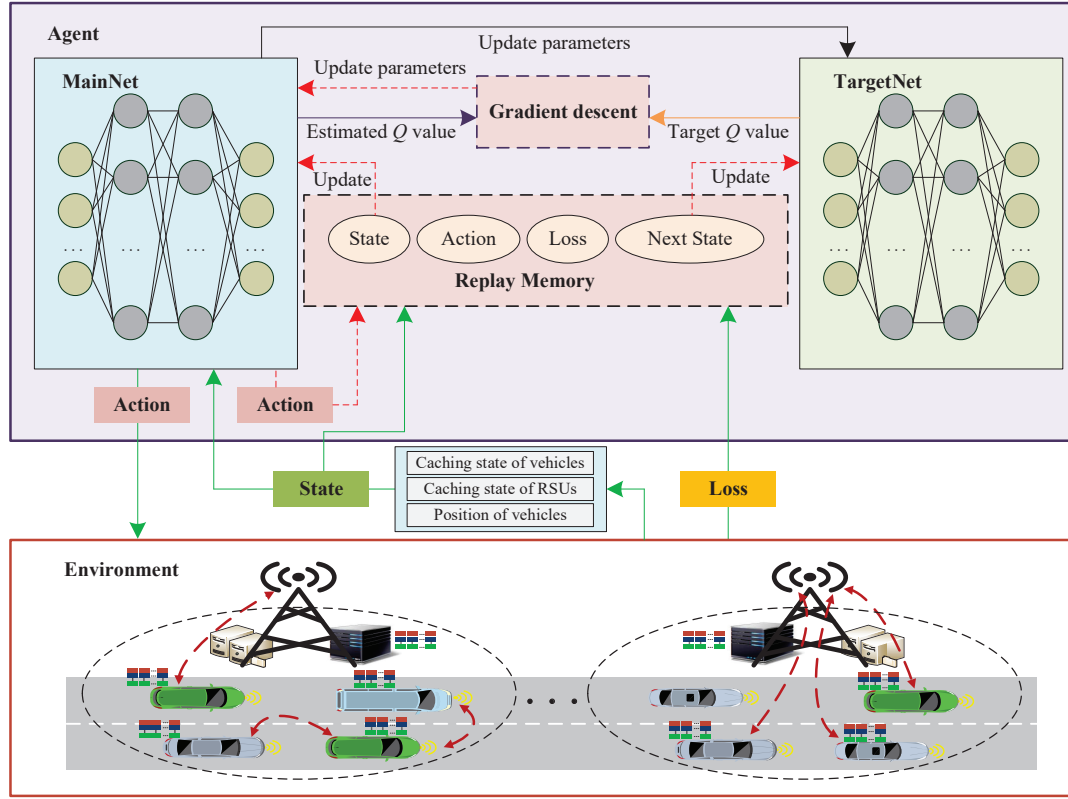


Fig. 3: DQN architecture for solving the data scheduling problem.

better data scheduling strategies, otherwise the action that has the minimum Q -value is chosen.

It is noteworthy that the same parameters are used for calculating the estimated and target Q -values. As a consequence, there is a big correlation between the estimated and target Q -values. Therefore, it means that at every step of training, our estimated Q -value shifts but also the target Q -value shifts, which lead to a big oscillation in training. To address this issue, we introduce a separate target Q -network, to calculate the target Q -value. The parameters of the target Q -network at time-slot t is denoted as $\bar{\theta}^t$. The target Q -value is correspondingly expressed as

$$\bar{y}^t = Loss^t + \eta \min_{A^{t+1}} Q(S^{t+1}, A^{t+1}; \bar{\theta}^t). \quad (31)$$

Similarly, y^t in formulas (28) and (29) is substituted by \bar{y}^t . And $\bar{\theta}^t$ is hold fixed and only updated with the DQN parameters (θ^t) every ζ time-slots. Fig. 3 shows the DQN architecture for solving the data scheduling problem in the VEC network, where MainNet refers to the neural network used to calculate the estimated Q -value and TargetNet refers to the neural network used to calculate the target Q -value. The full algorithm for training the DQN for optimal task scheduling is presented in Algorithm 1.

VI. SIMULATION RESULTS AND DISCUSSIONS

In this section, we conduct simulations to validate the performance of the proposed data scheduling scheme. Firstly, we describe the simulation scenario and parameter settings. Next, we discuss the simulation results.

Algorithm 1 DQN-Based Data Scheduling

- 1: Initialize replay memory \mathcal{D}
- 2: Initialize DQN with random weights θ
- 3: Initialize target DQN with weights $\bar{\theta} = \theta$
- 4: **for** episode $e = 1, \dots, e_{\max}$ **do**
- 5: Observe the initial state S^1 ;
- 6: **for** time-slot $t = 1, \dots, t_{\max}$ **do**
- 7: Choose a random probability p ;
- 8: **if** $p \leq \epsilon$ **then**
- 9: Select a random action A^t ,
- 10: **else**
- 11: Choose action $A^t = \arg \min_{A^t} Q(S^t, A^t; \theta^t)$;
- 12: **end if**
- 13: Execute action A^t , calculate $Loss^t$ and derive the next state S^{t+1} according to formulas (17)-(20);
- 14: Store the experience $(S^t, A^t, Loss^t, S^{t+1})$ into \mathcal{D} in the replay memory;
- 15: Get random minibatch of samples $(S^i, A^i, Loss^i, S^{i+1})$ from \mathcal{D} ;
- 16: Calculate the target Q -value from the target DQN, $\bar{y}^i = Loss^i + \eta \min_{A^{i+1}} Q(S^{i+1}, A^{i+1}; \bar{\theta}^i)$;
- 17: Perform the gradient descent step on $Er(\theta^i) = E[(\bar{y}^i - Q(S^i, A^i; \theta^i))^2]$ with respect to θ^i , and update θ^i ;
- 18: Every ζ time-slots, update $\bar{\theta}^t$ with θ^t ;
- 19: **end for**
- 20: **end for**

A. Simulation Setup

We consider a two-way three-lane scenario. The length of per lane is 1000 m and the width of per lane is 4 m. And one RSU is deployed in the middle of roadside, with coordinate (0 m, 0 m). Three vehicles that are on three different lanes are keeping moving back and forth along their lanes. The initial positions of the three vehicles are set to (-500 m, 2 m), (0 m, 6 m), and (500 m, 10 m), respectively. For the speed of the three vehicles, we use part of the GAIA Open Dataset containing speeds of DiDi Express in Xi'an China [33]. The data set contains the GPS coordinates and real-time speeds of DiDi Express over 30 days and over thousands of districts and roads. We randomly choose three loads and calculate the average speeds of vehicles over the three roads, respectively. Based on the speed statistics, we set the speeds of the three vehicles to 17.7 km/h, 35.8 km/h, and 52.6 km/h, respectively. The coverage radius of RSU and vehicles are set to 500 m and 250 m, respectively. The generated data of the vehicles is classified into four types. The length of time-slot is set to 100 ms. The data size is randomly distributed between 0.2 and 5.0, and the data delay constraints are randomly generated from {1, 2, 3, 4}. The detailed parameters setting about vehicles and RSU is shown in Table II.

TABLE II: Parameters setting about vehicles and RSU

Description	Parameter	Value
Coverage radius of vehicles	R^V	250 m
Coverage radius of RSU	R_k	500 m
Length of time-slot	Δt	100 ms
Number of data types	M	4
Amount of type- i data	D_i	0.2 ~ 5.0
Delay constraints of type- i data	T_i	{1, 2, 3, 4}
Bandwidth per channel	B	1 MHz
Number of V2I channels	C_{V2I}	20
Number of V2V channels	C_{V2V}	10
Transmission power of vehicles	P_n^{tr}	30 dBm
White Gaussian noise power	ω_0	-100 dBm
Path loss exponent	ϑ	4
Switched capacitance coefficient	κ_1, κ_2	$10^{-24}, 10^{-29}$
Processing capability of RSU	f_k^{ECN}	100 G cycles/s
Processing capability of vehicle	f_n^{local}	1.4 G cycles/s
Processing density of data	c	100 cycles/bit
Penalty coefficient	ξ	5
Energy consumption cost coefficient	ϱ	2.44×10^{-4}
V2I cost coefficient	c^{I}	1
V2V cost coefficient	c^{V}	0.5
Cost for RSU processing data	c^{ECN}	2

Hereinafter, we introduce the simulation settings about DQN. The simulation uses TensorFlow [34] to implement the DQN. We use a GPU-based server with 4 NVIDIA GTX 2080 Ti GPUs, where the CPU is Intel Xeno(R) E5-2690 v4 with 64G memory. Software environment we utilize is Keras 2.3.1 (based on TensorFlow) with Python 3.7 on Ubuntu 16.04.6 LTS. Both the main and target Q -networks use fully connected deep neural network with two hidden layers. The ReLU function is adopted as the activation function for hidden layers. The number of nodes for the hidden layers is set to 128. We set the learning rate, discount factor to 0.001 and 0.95, respectively. For the ϵ -greedy policy, we initialize $\epsilon = 1.0$ and let it decrease by a decay coefficient 0.995 over time-slots until it reaches 0.1. The minibatch and maximum replay memory

sizes are set to 64 and 10000, respectively. The maximum time-slots t_{\max} and episodes e_{\max} are set to 600 and 3000, respectively. The parameters updating frequency for target DQN ζ is set to 10. The detailed parameters setting of DQN is listed in Table III.

B. Simulation Results

a) Effectiveness: We mainly consider the following strategies, and evaluate their performance under the same conditions.

- Our proposal, that is, the proposed data scheduling scheme aiming at minimizing the cumulative value of loss over time-slots, and it is solved by DQN with separated target Q -network (recorded as *our proposal*);
- The proposed data scheduling scheme aiming at minimizing the cumulative value of loss over time-slots, and it is solved by DQN without separate target Q -network (recorded as *DQN without separate target Q-network*).

We first evaluate the timeliness of *our proposal*. Similarly, the scheduling strategy *without separate target Q-network* is evaluated under the same conditions as comparison scheme. Fig. 4 shows the changes in the reward (i.e., the data scheduling loss) obtained by the deep reinforcement learning agent and training episodes under different schemes. The figure shows that the loss of *our proposal* tends to be optimal and stable in about 200 episodes of training, and no divergence and oscillation problems are encountered in the simulation. As the episodes increase, the loss gradually stabilizes at the optimal value, which means that the agent of *our proposal* has learned the optimal data scheduling strategy to minimize long-term loss. The scheme *DQN without separate target Q-network* requires longer episodes to stabilize the loss (approximately 400). Although being stable, the loss of the scheme *DQN without separate target Q-network* is prone to relatively fluctuations and oscillation. This is because there is a big correlation between the estimated and target Q -values, and the estimated Q -values shift but also the target Q -values shifts when the same parameters are used during training process. The figure indicates that *our proposal* is more effective for solving the data scheduling problem.

b) Average Loss: Average loss indicates the average value of long-term data scheduling loss (i.e., *Loss* in formula (11)) over time-slots. The following strategies are used for performance comparison:

- The proposed data scheduling scheme based on DQN with separate target Q -network (Record as *our proposal*);
- The data is only computed and processed by vehicles themselves (Record as *local-pro-only*);
- The data is only offloaded to RSU for computing and processing (Record as *offload-only*);
- The scheme similar to our proposal, except for the collaborative computing when V2V action is adopted. That is to say, this scheme has no collaborative computing (Record as *without collaborative computing*).

We evaluate the performance of average loss under different schemes and different data size. The data size indicates the amount of data generated in each time-slot. The overall results

TABLE III: Parameters setting of DQN

Description	Parameter	Value
State vector	S^t	$\{Q_1^t, Q_2^t, Q_3^t, G^t, \Phi^t\}$
Action vector	A^t	$\{a_1^t, a_2^t, a_3^t, \bar{a}^t\}$
Network architecture		Fully connected neural network
Number of hidden layers		2
Number of nodes for the first hidden layer		128
Number of nodes for the second hidden layer		128
Activation function for hidden layers		ReLU
Learning rate		0.001
Discount factor	η	0.95
Initial ϵ		1
ϵ decay coefficient		0.995
Minimal ϵ		0.1
Minibatch size		64
Maximum replay memory size	$ \mathcal{D} $	10000
maximum time-slots	t_{\max}	600
Maximum episodes	e_{\max}	3000
Parameter updating frequency for target DQN	ζ	10

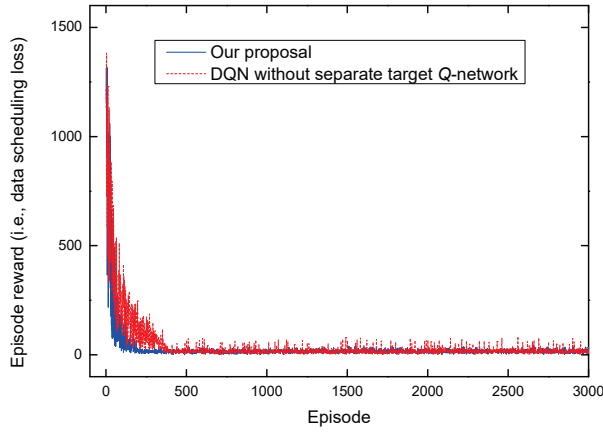
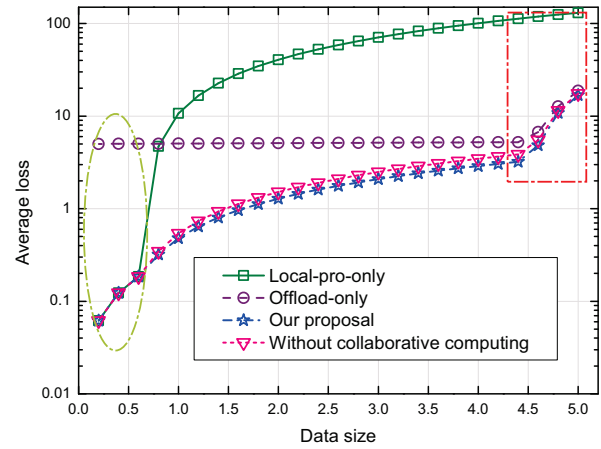


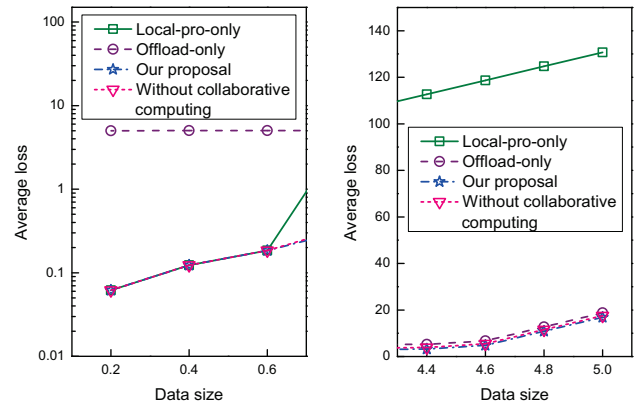
Fig. 4: Episode reward (i.e., data scheduling loss) achieved during training.



(a) Data size varied from 0.2 to 5.0.

are shown in Fig. 5. Overall, the average loss increases with the increasing data size. Specifically, the variation process of average loss can be divided into four phases. The first phase is when data size varies from about 0.2 to 0.6. The second phase is when data size varies from about 0.6 to 1.0. The third phase is when data size varies from about 1.0 to 4.4. The fourth phase is when data size varies from 4.4 to 5.0. In the following, the four phases will be discussed in detail, respectively.

The first phase is shown as the circled part in Fig. 5(a). To show the variation process clearly, this part is zoomed in, as shown in the left figure of Fig. 5(b). In this phase, the average losses of *local-pro-only*, *our proposal*, and *without collaborative computing* have the same value and variation trending. However, the average loss of *offload-only* is much higher than other three schemes. This is because the local computing capacity of vehicles is sufficient to process the data when data size is small, hence, there will be no penalty caused by unprocessed data. And since the cost of local computing is much lower than the cost of RSU computing (consisted of V2I communication cost and RSU computing cost), both agents



(b) The circled part and framed part in the Fig. 5(a).

Fig. 5: Average loss with various data size and different schemes.

of *our proposal* and *without collaborative computing* learn to adopt local computing actions for the purpose of minimizing the cost thus minimizing the loss. The underlying reason why the average loss of *offload-only* increases slowly is that the

cost of V2I communication and RSU computing is relatively high. In each time-slot, the data of vehicles is offloaded to RSU through V2I communication, and RSU computes the offloaded data. The communication and computation cost are relatively fixed because all the costs are calculated in the unit of time-slot. And the cost of computing energy consumption of RSU is much lower than the communication and computation costs. Accordingly, the increase of data size would result in small increase in loss caused by the increase of computing energy consumption. Another reason why the average loss looks like it increases slowly is that the vertical scale is of exponential type. We tabulate the variation of average loss of *offload-only* in Table IV.

TABLE IV: The average loss of *offload-only* under small data sizes

Data size	0.2	0.4	0.6	0.8
Average loss	5.0051	5.0158	5.0268	5.0379

For the second phase, as the data size increases, the average loss of *local-pro-only* increases rapidly and exceeds that of *offload-only*. The underlying reason is that the data size in this phase exceeds the local computing capacity of vehicles, the task not meeting the delay constraints will result in a high loss because of the penalty to the system, as described in formulas (10) and (11). It is noteworthy that the average losses of both *our proposal* and *without collaborative computing* are lower than that of *local-pro-only* and *offload-only*. This is because the former two schemes can optimally decide whether data is processed locally or offloaded to RSU, hence the data that exceeds the local computing capacity is offloaded to RSU, for the purpose of minimizing the loss. In addition, the average loss of *our proposal* is lower than that of *without collaborative computing*. The underlying reason is that, in *our proposal*, the data can be migrated to collaborative vehicles through V2V communication for a lower cost thus a lower loss than be offloaded to RSU.

As data size keeps increasing, the third phase depicts the increase of average loss for all the four schemes. In this phase, the data size does not exceed the sum of local computing capacity and the data transmission capacity. In *our proposal* and *without collaborative computing*, the data can be partly offloaded to RSU for processing. Accordingly, no penalty is introduced, which has a lower average loss than *local-pro-only* and *offload-only*.

The fourth phase is shown as the framed part in Fig. 5(a). In this phase, the average losses of *offload-only*, *our proposal*, and *without collaborative computing* have a rapid increase when the data size increases. The underlying reason is that, during this phase, the data size exceeds the sum of local computing capacity and the data transmission capacity. In this case, penalty will be introduced due to unprocessed data, which results in a high loss according to formula (11). To show the variation process clearly, this part is also zoomed in, and the Y-axis is changed to linear scale, as shown in the right figure of Fig. 5(b). When the data size in this phase reach a threshold (approximately 4.6), the average losses of *offload-only*, *our proposal*, and *without collaborative computing* have the same

increasing trending as *local-pro-only*. This is because all the four schemes can not completely process all the generated data when the data size exceed a threshold. Accordingly, how fast the average loss increases depends on how fast the data size increases. Furthermore, we can find that the average loss of *our proposal* has almost the same values as that of the scheme without collaborative computing. This is because if the data is migrated to collaborative vehicles, high loss will be caused due to the limited computing capacity of collaborative vehicles. Instead, during training process, the agent learns that the data should be offloaded to RSU for processing to minimize the loss.

c) Data scheduling Actions: To understand how the agent of our proposal selects the optimal data scheduling actions during training process, we evaluate the proportion of scheduling actions on both vehicle side and RSU side. As comparison, we select four different data sizes, i.e., 0.5, 2.0, 3.5 and 5.0.

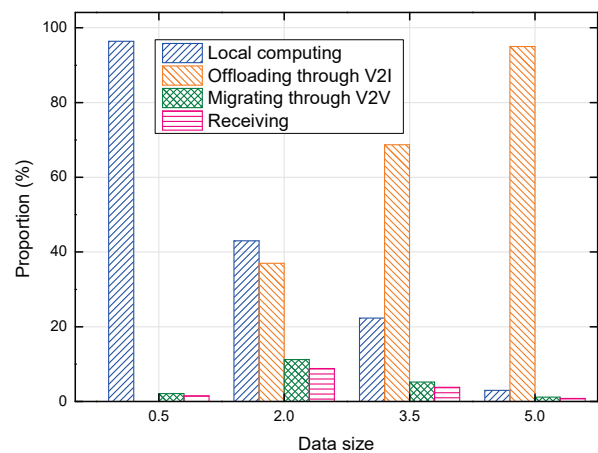


Fig. 6: Average proportion of data scheduling actions on vehicle side under various data size.

Fig. 6 shows the proportion of data scheduling actions on vehicle side under different data size. It is not difficult to find out that the action is gradually changing from local computing to offloading through V2I as data size increases. The reason is that when data size is small, the local computing capacity of vehicles is sufficient to process the data, the agent prefers to select local computing action for vehicles to minimize the cost thus minimize the loss. As the data size exceeds the local computing capacity of vehicles, the proportions of offloading through V2I, migrating through V2V, and receiving actions increase. This is because the task date should be processed under delay constraints to reduce the timeout penalty thus minimize the loss. And the proportion of migrating through V2V action has almost the same value as the receiving action. The reason is that the establishment of a communication link must contain a sender and a receiver. It is noteworthy that the proportion of migrating through V2V action increases first then decreases as data size increases. The underlying reason is that, when data size is small, local computing is sufficient and the cost of migrating through V2V is higher than local computing. As data size increases and exceeds the local

computing capacity, instead of offloading to RSU, data can be migrated to collaborative vehicles for processing because the cost of migrating through V2V is lower than offloading through V2I. As data size keeps increasing, the proportion of migrating through V2V decreases until its value reaches almost 0. This is because if migrating through V2V action is adopted, the limited computing capacity of collaborative vehicle will cause much timeout penalty thus greatly increasing the loss, which violates agent's purpose during training process.

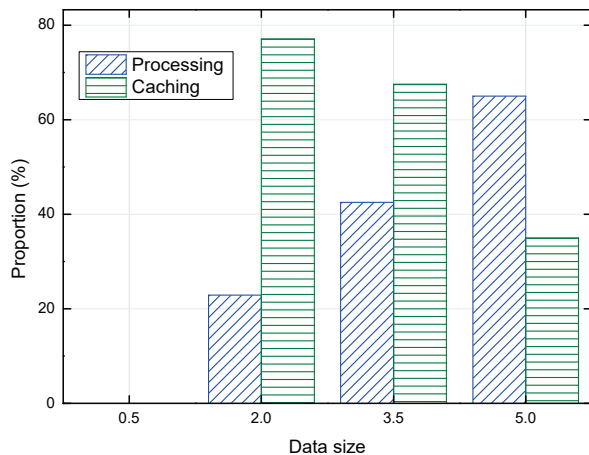


Fig. 7: Average proportion of data scheduling actions on RSU side under various data size.

Fig. 7 shows the proportion of data scheduling actions on RSU side under different data size. There are two actions (i.e., processing data and caching data) related to RSU. It can be seen that the proportions of both processing and caching equal 0 when the data size is 0.5. The underlying reason is that under this value of data size, the local computing capacity of vehicles is sufficient to process data, the agent learns to choose local computing action during training to reduce cost thus minimize the loss. It is also consistent with the analysis for the first phase in Fig. 5. More data is offloaded to RSU as data size increases, the proportion of processing action also increases. This is because that there is a high probability that more data is cached in the queue with index 1 as the data size increases, penalty would be thus introduced if data is not processed. In order to reduce penalty thus minimize the loss, the proportion of processing data increases with the increase of data size. We can also draw a valuable conclusion that even when the data size is 5.0, the computing capacity of RSU is sufficient to process data. The underlying reason is that the proportion of processing action is not 100% when data size is 5.0, which indicates that the RSU processes data not in each time-slot to reduce the cost of processing data thus minimize the loss.

VII. CONCLUSION

We argue that with the development of edge computing and autonomous driving, data communication and computation offloading will equally dominate the usage of bandwidth, and a

unified framework to jointly consider all resources available in the network is necessary. As motivated, we have investigated the data scheduling problem in a unified paradigm where data can be processed locally, offloaded to RSUs, migrated to collaborative vehicles, or kept in caching queues. A multi-queue model for data caching on both vehicle side and RSU side has been developed accordingly. To derive the optimal data scheduling strategy, we formulated a MDP model to reflect the interaction between data scheduling actions and loss, by jointly considering communication resource, caching state, computing resource, mobility of vehicles as well as delay requirements of data. Next, we developed a DRL framework for the data scheduling problem, refining the state, action, and reward function, respectively. A DQN-based method with separate target Q -network was then proposed for solving this problem, which can efficiently optimize the data scheduling for minimizing the loss. Extensive simulation results were presented to illustrate that our proposal efficiently reduces data processing cost and helps data be processed under delay constraints.

In our future work, we intend to extend in three directions. First, we would include the migration of data between RSUs and consider the edge-cloud collaboration in the vehicular networks. Second, we would apply the proposed algorithm in a real-world testbed and verify the performance. Third, since overhead would be generated when the state of each vehicle is collected, we would discuss and evaluate the overhead performance of the proposed method.

REFERENCES

- [1] L. Liu, C. Chen, Q. Pei, S. Maharjan, and Y. Zhang, "Vehicular edge computing and networking: A survey," *arXiv: 1908.06849*, 2019. [Online]. Available: <https://arxiv.org/abs/1908.06849>
- [2] K. Zhang, S. Leng, X. Peng, L. Pan, S. Maharjan, and Y. Zhang, "Artificial intelligence inspired transmission scheduling in cognitive vehicular communications and networks," *IEEE Internet Things J.*, vol. 6, no. 2, pp. 1987-1997, 2019.
- [3] Y. Dai, D. Xu, S. Maharjan, G. Qiao, and Y. Zhang, "Artificial intelligence empowered edge computing and caching for internet of vehicles," *IEEE Wireless Commun.*, vol. 26, no. 3, pp. 12-18, 2019.
- [4] P. Mach, and Z. Becvar, "Mobile edge computing: A survey on architecture and computation offloading," *IEEE Commun. Surv. Tutor.*, vol. 19, no. 3, pp. 1628-1656, 2017.
- [5] C. Li, Q. Luo, G. Mao, M. Sheng, and J. Li, "Vehicle-mounted base station for connected and autonomous vehicles: opportunities and challenges," *IEEE Wireless Commun.*, vol. 26, no. 4, pp. 30-36, 2019.
- [6] F. Sun, F. Hou, N. Cheng, M. Wang, H. Zhou, L. Gui, and X. Shen, "Cooperative task scheduling for computation offloading in vehicular cloud," *IEEE Trans. Veh. Technol.*, vol. 67, no. 11, pp. 11049-11061, 2018.
- [7] C. Wang, C. Liang, F. R. Yu, Q. Chen, and L. Tang, "Computation offloading and resource allocation in wireless cellular networks with mobile edge computing," *IEEE Trans. Wireless Commun.*, vol. 16, no. 8, pp. 4924-4938, 2017.
- [8] K. Zhang, Y. Mao, S. Leng, S. Maharjan, and Y. Zhang, "Optimal delay constrained offloading for vehicular edge computing networks," in *Proc. IEEE Int. Conf. Commun. (IEEE ICC)*, 2017, pp. 1-6.
- [9] P. Paymard, S. Rezvani, and N. Mokari, "Joint task scheduling and uplink/downlink radio resource allocation in PD-NOMA based mobile edge computing networks," *Phys. Commun.*, vol. 32, no. 2019, pp. 160-171, 2018.
- [10] N. Cheng, F. Lyu, W. Quan, C. Zhou, H. He, W. Shi, and X. Shen, "Space/Aerial-assisted computing offloading for iot applications: A learning-based approach," *IEEE JSAC*, vol. 37, no. 5, pp. 1117C1129, 2019.

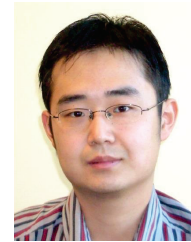
- [11] Q. Qi, and Z. Ma, "Vehicular edge computing via deep reinforcement learning," *arXiv: 1901.04290*, 2018. [Online]. Available: <https://arxiv.org/abs/1901.04290>
- [12] J. Chen, and X. Ran, "Deep learning with edge computing: A review," *Proc. IEEE*, vol. 107, no. 8, pp. 1655-1674, 2019.
- [13] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge intelligence: Paving the last mile of artificial intelligence with edge computing," *Proc. IEEE*, vol. 107, no. 8, pp. 1738-1762, 2019.
- [14] L. T. Tan, and R. Q. Hu, "Mobility-aware edge caching and computing in vehicle networks: A deep reinforcement learning," *IEEE Trans. Veh. Technol.*, vol. 67, no. 11, pp. 10190-10203, 2018.
- [15] Y. He, N. Zhao, and H. Yin, "Integrated networking, caching, and computing for connected vehicles: A deep reinforcement learning approach," *IEEE Trans. Veh. Technol.*, vol. 67, no. 1, pp. 44-55, 2018.
- [16] K. Zhang, Y. Zhu, S. Leng, Y. He, S. Maharjan, and Y. Zhang, "Deep learning empowered task offloading for mobile edge computing in urban informatics," *IEEE Internet of Things J.*, vol. 6, no. 5, pp. 7635-7647, 2019.
- [17] K. Zhang, Y. Mao, S. Leng, Y. He, and Y. Zhang, "Mobile-edge computing for vehicular networks: A promising network paradigm with predictive off-loading," *IEEE Veh. Technol. Mag.*, vol. 12, no. 2, pp. 36-44, 2017.
- [18] T. Q. Dinh, J. Tang, Q. D. La, and T. Q. S. Quek, "Offloading in mobile edge computing: Task allocation and computational frequency scaling," *IEEE Trans. Commun.*, vol. 65, no. 8, pp. 3571-3584, 2017.
- [19] J. Du, F. R. Yu, X. Chu, J. Feng, and G. Lu, "Computation offloading and resource allocation in vehicular networks based on dual-side cost minimization," *IEEE Trans. Veh. Technol.*, vol. 68, no. 2, pp. 1079-1092, 2019.
- [20] Q. Luo, C. Li, T. H. Luan, and Y. Wen, "Optimal utility of vehicles in LTE-V scenario: An immune clone-based spectrum allocation approach," *IEEE Trans. Intell. Transp. Sys.*, vol. 20, no. 5, pp. 1942-1953, 2018.
- [21] M.-A. Messous, H. Sedjelmaci, N. Houari, and S.-M. Senouci, "Computation offloading game for an UAV network in mobile edge computing," in *Proc. IEEE Int. Conf. Commun. (IEEE ICC)*, 2017, pp. 1-6.
- [22] T. Q. Dinh, Q. D. La, T. Q. S. Quek, and H. Shin, "Learning for computation offloading in mobile edge computing," *IEEE Trans. Commun.*, vol. 66, no. 12, pp. 6353-6367, 2018.
- [23] L. Lin, X. Liao, H. Jin, and P. Li, "Computation offloading toward edge computing," *Proc. IEEE*, vol. 107, no. 8, pp. 1584-1607, 2019.
- [24] J. Li, G. Luo, N. Cheng, Q. Yuan, Z. Wu, S. Gao, and Z. Liu, "An end-to-end load balancer based on deep learning for vehicular network traffic control," *IEEE Internet Things J.*, vol. 6, no. 1, pp. 953-966, 2019.
- [25] E. Li, L. Zeng, Z. Zhou, X. Chen, "Edge AI: on-demand accelerating deep neural network inference via edge computing," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 447-457, 2019.
- [26] H. Cao, J. Cai, "Distributed multiuser computation offloading for cloudlet-based mobile cloud computing: A game-theoretic machine learning approach," *IEEE Trans. Veh. Technol.*, vol. 67, no. 1, pp. 752-764, 2017.
- [27] Y. Mao, J. Zhang, S. H. Song, and K. B. Letaief, "Power-delay tradeoff in multi-user mobile-edge computing systems," in *Proc. IEEE Global Communications Conference (GLOBECOM)*, 2016.
- [28] Y. Wang, M. Sheng, X. Wang, L. Wang, and J. Li, "Mobile-edge computing: Partial computation offloading using dynamic voltage scaling," *IEEE Trans. Commun.*, vol. 64, no. 10, pp. 4268-4282, 2016.
- [29] Y. Kim, J. Kwak, and S. Chong, "Dual-side optimization for cost-delay tradeoff in mobile edge computing," *IEEE Trans. Veh. Technol.*, vol. 67, no. 2, pp. 1765-1781, 2018.
- [30] O.-A. Maillard, T. A. Mann, and S. Mannor, "How hard is my MDP? The distribution-norm to the rescue," in *Proc. Adv. Neural Inf. Process. Syst.*, pp. 1835-1843, 2014.
- [31] V. Mnih, K. Kavukcuoglu, D. Silver, *et al.*, "Playing atari with deep reinforcement learning," *arXiv: 1312.5602*, 2013. [Online]. Available: <https://arxiv.org/abs/1312.5602>
- [32] V. Mnih, K. Kavukcuoglu, D. Silver, *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529-533, 2015.
- [33] Didi Chuxing GAIA Initiative, "Urban traffic time index and trajectory data (new)," Downloaded from <https://gaia.didichuxing.com>.
- [34] M. Abadi *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous systems", *arXiv:1603.04467*, 2016. [Online]. Available: <https://arxiv.org/abs/1603.04467>



in vehicular networks.



mobile ad hoc networks, and wireless sensor networks.



the content distribution in vehicular networks, mobile cloud computing, and fog computing.



Dr. Shi received the Most Downloaded Paper Award in IEEE Computer. He is a recipient of the National Outstanding Ph.D. dissertation award of China (2002), the NSF CAREER award (2007), Wayne State University Career Development Chair award (2009), Charles H. Gershenson Distinguished Faculty Fellow (2015), College of Engineering Faculty Research Excellence Award (2016), the Best Paper award of ICWE'04, IEEE IPDPS'05, HPCChina'12, IEEE IISWC'12, IEEE eHealth'17, the Best Paper Nominee award of ACM UbiComp'14, and the Best Student Paper Award of IEEE HealthCom'15. He is the Editor-in-Chief of Elsevier Smart Health, the associate Editor-in-Chief of IEEE Internet Computing, the associate editor of IEEE Transactions on Services Computing, ACM Transactions on Internet of Things, Springer Computing Journal, the editor of IEEE Blockchain Newsletter, IEEE Internet Computing, SUSCOM, JCST, IJNet. He served as the Program/General Chair/Co-Chair for the MetroCAD'19, MetroCAD'18, ACM/IEEE SEC'18, ACM/IEEE CHASE'18, ACM/IEEE CHASE'16, IEEE HotWeb'15, IEEE CSE'11. He is an IEEE Fellow and a Distinguished Scientist of ACM.

Qu Yuan Luo received the B.E. degree from Changchun University of Science and Technology, Changchun, China, in 2011. He is currently working toward the Ph.D. degree in communication and information system with Xidian University, Xi'an, China. He is also a Visiting Scholar with Computer Science, Wayne State University, USA, from 2019. His collaborator at Wayne State University is Prof. Weisong Shi. His current research interests include intelligent transportation systems, content distribution, edge computing and resource allocation

Changle Li (M'09-SM'16) received the Ph.D. degree in communication and information system from Xidian University, Xi'an, China, in 2005. He conducted his postdoctoral research in Canada and the National Institute of information and Communications Technology, Japan, respectively. He had been a Visiting Scholar with the University of Technology Sydney and is currently a Professor with the State Key Laboratory of Integrated Services Networks, Xidian University. His research interests include intelligent transportation systems, vehicular networks,

Tom H. Luan (M'10) received the B.E. degree from Xi'an Jiaotong University, Xi'an, China, in 2004; the M.Phil. degree from the Hong Kong University of Science and Technology, Kowloon, Hong Kong, in 2007; and the Ph.D. degree from the University of Waterloo, Waterloo, ON, Canada, in 2012. He has been a Lecture with the School of Information Technology, Deakin University, Burwood, VIC., Australia from 2013 to 2017 and is currently a Professor with the School of Cyber Engineering, Xidian University. His current research interests focus on

Weisong Shi (F'16) received the B.S. degree from Xidian University, Xi'an, China, in 1995, and the Ph.D. degree from the Chinese Academy of Sciences, in 2000, both in computer engineering. He is a Professor of Computer Science with Wayne State University, where he is the associate dean for research of College of Engineering. His research interests include edge computing, big data systems, computer systems, energy-efficiency computer systems, mobile and connected health, connected and autonomous driving.