

Distributed Deep Multi-Agent Reinforcement Learning for Cooperative Edge Caching in Internet-of-Vehicles

Huan Zhou[✉], Member, IEEE, Kai Jiang[✉], Student Member, IEEE, Shibo He[✉], Senior Member, IEEE, Geyong Min[✉], Senior Member, IEEE, and Jie Wu[✉], Fellow, IEEE

Abstract—Edge caching is a promising approach to reduce duplicate content transmission in Internet-of-Vehicles (IoVs). Several Reinforcement Learning (RL) based edge caching methods have been proposed to improve the resource utilization and reduce the backhaul traffic load. However, they only obtain the local sub-optimal solution, as they neglect the influence from environments by other agents. This paper investigates the edge caching strategies with consideration of the content delivery and cache replacement by exploiting the distributed Multi-Agent Reinforcement Learning (MARL). A hierarchical edge caching architecture for IoVs is proposed and the corresponding problem is formulated with the goal to minimize the long-term content access cost in the system. Then, we extend the Markov Decision Process (MDP) in the single agent RL to the context of a multi-agent system, and tackle the corresponding combinatorial multi-armed bandit problem based on the framework of a stochastic game. Specifically, we firstly propose a Distributed MARL-based Edge caching method (DMRE), where each agent can adaptively learn its best behaviour in conjunction with other agents for intelligent caching. Meanwhile, we attempt to reduce the computation complexity of DMRE by parameter approximation, which legitimately simplifies the training targets. However, DMRE is enabled to represent and update the parameter by creating a lookup table, essentially a tabular-based method, which generally performs inefficiently in large-scale scenarios. To circumvent the issue and make more expressive parametric models, we incorporate the technical advantage of the Deep-Q Network into DMRE, and further develop a computationally efficient method (DeepDMRE) with neural network-based Nash

equilibria approximation. Extensive simulations are conducted to verify the effectiveness of the proposed methods. Especially, DeepDMRE outperforms DMRE, *Q*-learning, LFU, and LRU, and the edge hit rate is improved by roughly 5%, 19%, 40%, and 35%, respectively, when the cache capacity reaches 1,000 MB.

Index Terms—Edge caching, Internet-of-Vehicles, content delivery, cache replacement, multi-agent reinforcement learning.

I. INTRODUCTION

WITH the explosive growth of vehicles on the road and the progress of wireless multi-access technology, we have witnessed unprecedented changes in the traditional transportation system, which has evolved from a technology-driven era to a more powerful data-driven intelligent era [1]. As a fundamental paradigm of 5G networks, Internet-of-Vehicles (IoVs) enables a wide variety of vehicles to provide reliable vehicular multimedia services, cooperative cruise control, and path navigation, etc., which have paved a path towards intelligent transportation, and improved the driving safety and travel comfort of passengers [2], [3].

Meanwhile, these flourishing vehicular applications require vehicles to access vast amounts of Internet data, especially for some delay-sensitive contents, leading to severe network congestion and considerable delay in content delivery [4], [5], [6]. Technically, tight Quality-of-Service (QoS) requirements are generally placed for such applications, which makes cloud-based processing architectures infeasible due to the long transmission distance and limited backhaul link capacity inherent in content delivery from remote cloud servers. Indeed, the rapid increase in delay and unreliability driven by the transmission distance has become an inevitable issue for supporting massive content delivery and gained widespread attention lately [7], [8]. Despite the continuous efforts made on improving the backhaul link capacity through advanced technologies, the radio spectrum utilization has obviously reached the theoretical bound [9], [10], [11]. Therefore, these efforts are insufficient to cope with these great challenges. A fundamental innovation that breaks through the bottleneck of massive content delivery in IoVs there is urgently required.

Fortunately, large-scale data analysis shows that different contents often require different priorities. Only a few popular contents account for the majority of downloads, while the access demand for most of the rest is relatively

Manuscript received 15 January 2022; revised 18 July 2022, 28 October 2022, and 1 March 2023; accepted 17 April 2023. This work was supported in part by the National Natural Science Foundation of China under Grant 62172255 and Grant U1909207 and in part by the Outstanding Youth Program of Hubei Natural Science Foundation under Grant 2022CFA080. The associate editor coordinating the review of this article and approving it for publication was A. Schmeink. (*Corresponding author: Kai Jiang*.)

Huan Zhou is with the School of Computer Science, Northwestern Polytechnical University, Xi'an 710129, China, and also with the College of Computer and Information Technology, China Three Gorges University, Yichang 443002, China (e-mail: zhouchuan17@gmail.com).

Kai Jiang is with the School of Cyber Science and Engineering, Wuhan University, Wuhan 430000, China (e-mail: kai.jiang@whu.edu.cn).

Shibo He is with the College of Control Science and Technology, Zhejiang University, Hangzhou 310027, China (e-mail: s18he@zju.edu.cn).

Geyong Min is with the Department of Computer Science, College of Engineering, Mathematics, and Physical Sciences, University of Exeter, EX4 4QF Exeter, U.K. (e-mail: g.min@exeter.ac.uk).

Jie Wu is with the Department of Computer and Information Sciences, Temple University, Philadelphia, PA 19122 USA (e-mail: jiewu@temple.edu).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TWC.2023.3272348>.

Digital Object Identifier 10.1109/TWC.2023.3272348

72 small [12], [13]. This request pattern promotes the imple-
 73 mentation of edge caching technology in IoVs. Expressly,
 74 edge caching has provided an alternative to alleviate the
 75 backhaul link strain and content access delay inside the
 76 future IoVs, which pre-caches frequently-used contents close
 77 to vehicles by pushing cloud functions to intermediate Road-
 78 side Units (RSUs). Hence, edge caching enables vehicles to
 79 access popular content from the caching-enabled nearby RSUs
 80 directly, instead of repeatedly downloading from remote cloud
 81 servers [14]. In this way, the redundant traffic and transmission
 82 resource consumption at both backhaul and core networks can
 83 be significantly reduced, and the QoS is improved as well.
 84 In addition, enabling edge caching in IoVs can benefit extra
 85 technical superiorities from its distributed architectures and
 86 small-scale nature, including privacy protection, scalability,
 87 context awareness, and so on.

88 Compared with the cloud server, the cache capacity of a
 89 single edge node is insufficient to support all popular con-
 90 tents. Without edge cooperation, the overall cache utilization
 91 and effectiveness are prone to be under-utilized. Specifically,
 92 adjacent RSUs can share data traffic and exchange popular
 93 contents rather than operating individually, so the bilateral
 94 synergy between their caches and the centralized cloud should
 95 also be leveraged to facilitate good performance, which largely
 96 depends on a well-designed cooperative edge caching strategy.
 97 On the other hand, with the revival of artificial intelligence,
 98 Reinforcement Learning (RL) has exhibited its particular
 99 potential for the method design of efficient edge caching in
 100 IoVs. In RL, the agent can well capture the hidden dynamics
 101 of the environments and imitate the best features by trial-
 102 and-error interaction, thus learning a series of intelligence
 103 behaviors to enhance the edge cache utilization [15], [16], [17].
 104 However, massively diverse, highly vibrant and distributed
 105 edge caching contexts make most previous RL work unable
 106 to adapt as they neglected the environment's influence by
 107 other agents when an agent interacts and learns the settings
 108 independently. Therefore, only the local sub-optimal solution
 109 can be obtained in the system, not the optimal global one,
 110 especially in a long-term optimization problem [18].

111 Indeed, it is more reasonable and helpful to investigate from
 112 a distributed perspective under the multi-agent context, which
 113 leads to inherent robustness and high scalability.¹ Meanwhile,
 114 several studies [19], [20], and [21] have proved that agents
 115 considering the impact of joint actions perform better than
 116 corresponding agents learning only based on their own actions.
 117 Therefore, this paper investigates the edge caching strategy
 118 with consideration of the content delivery and cache replace-
 119 ment by exploiting the distributed Multi-Agent Reinforcement
 120 Learning (MARL). A hierarchical edge caching architecture
 121 for IoVs is first proposed, where the cooperative caching
 122 among multi-RSUs and Macro Base Station (MBS) is utilized
 123 to reduce the content access cost and the backhaul traffic in
 124 the system. Then, the Markov Decision Process (MDP) in
 125 the single-agent RL is extended to the context of multi-agent

¹Robustness: when one or more agents fail in the system, the remaining agents can take over some of their tasks.

Scalability: by sure design, most multi-agent systems allow easy insertion of new agents.

system and the corresponding combinatorial multi-armed bandit problem is tackled based on the framework of stochastic games. Specifically, a **Distributed MARL-based Edge caching method (DMRE)** is proposed firstly, where each agent can adaptively learn its best behaviour in conjunction with other agents for intelligent caching. Meanwhile, we attempt to reduce the computation complexity of DMRE by parameter approximation, which legitimately simplifies the training targets. However, DMRE is enabled to represent and update the parameter by creating a lookup table, essentially a tabular-based method, which generally performs inefficiently in large-scale scenarios. To circumvent the issue and make more expressive parametric models, we incorporate the technical advantage of the Deep-*Q* Network into DMRE, and further develop a computationally efficient method (**DeepDMRE**) with neural network-based Nash equilibria approximation. The main contributions are summarized as follows:

- 1) An original hierarchical edge caching architecture for IoVs is presented and the corresponding problem is formulated with the goal to minimize the long-term content access cost in the system.
- 2) To address the overhead minimization problem, an MDP for the optimization of cache replacement process in an available RSU is defined. Then, the MDP is extended to the context of a multi-agent system, and the corresponding combinatorial multi-armed bandit problem is tackled with the proposed DMRE firstly.
- 3) Furthermore, to circumvent the issue by DMRE and make more expressive parametric models, we incorporate the technical advantage of the Deep-*Q* Network into DMRE, and further develop a computationally efficient method (DeepDMRE) with neural network-based Nash equilibria approximation.
- 4) Extensive simulation results demonstrate that DeepDMRE significantly outperforms other benchmark methods in different scenarios. Meanwhile, relevant theoretical proofs of convergence and feasibility of the proposed methods are presented at the end of the paper.

The remainder of this paper is organized as follows. Section II and Section III describe the related work and the system framework, respectively. Section IV formulates the problem to minimize the long-term content access cost in the system. Section V extends MDP to a multi-agent system under the stochastic game framework, and proposes two MARL-based methods. Moreover, simulation results are analyzed in Section VI. Conclusion is summarized in Section VII.

II. RELATED WORK

Recently, extensive works have focused on content caching strategies at the edge of networks [22], [23], [24], [25], [26], [27], [28], [29], [30]. Jedari et al. [22] elaborated an exhaustive review on edge caching and discussed its implementation and outlook. Wu and Lu [23] investigated the tradeoff between content delivery delay and power consumption in small-cell network caching, and proposed an iterative algorithm to approach the optimal tradeoff between these two metrics. Yan et al. [24] exploited the energy assessment

models for mobile edge caching and proposed a predictive caching strategy to potentially improve the cache hit rates. Zhang et al. [25] focused on a learning-based edge caching scheme to enable cooperation among different edge nodes with limited resources, intending to reduce the content delivery latency and maximize the overall content caching value. Similarly, to maximize the utility of the caching system, a vehicular edge caching mechanism based on user preference similarity and service availability was proposed in [26]. Furthermore, Zhao et al. [27] dynamically orchestrated the cache resources in IoVs by leveraging the combined power of Lyapunov optimization, matching theory, and consensus alternating direction method of multipliers. In [28], Ao et al. leveraged the idea of coordinated multi-point and devised an efficient cooperative caching strategy to maximize the system throughput, where the available cache size and the network topology are both taken into account. In [29], Gu et al. proposed a dual-connectivity sub-6 GHz and mmWave heterogeneous network architecture to conduct a collaborative edge resource design, aiming to maximise virtual reality delivery reliability. In [30], Zhang et al. considered the stochastic geometry based probabilistic caching to improve the offloading rate for backhaul links, and presented an improved caching probability conversion algorithm to obtain the closed-form solutions. These works prove the advantage of edge caching. However, most of them are quasi-static and myopic, and thus perform unsatisfactorily in the diverse and dynamic IoVs system.

As Reinforcement Learning (RL) can well capture the hidden dynamics of the environments, it has been exploited to enhance the intelligence of IoVs. Specifically, Wang et al. [31] clustered vehicles to simplify the process of content requesting via the K-means method, and proposed an RL-based cooperative caching strategy with content request prediction in IoVs. In [32], Dai et al. developed a novel edge caching architecture for AI-empowered vehicular networks. They formulated a joint optimization of edge computing and caching to maximize the system utility, and exploited an RL based method to dynamically allocate computing and caching resources. Without any information or assumptions about the environment, Tan and Hu [33] achieved an RL based adaptive framework to tackle the resource allocation of communication, caching, and computing in vehicular networks. This framework was also used by Qiao et al. [34] to study the problem of content placement and content delivery in IoVs, which aims to reduce the system cost under the constraint of content delivery latency. Furthermore, Wang et al. [35] proposed a federated deep RL-based cooperative edge caching framework to eliminate duplicate traffic and improve the edge cache utilization with specific QoS requirements. Tian et al. [36] considered the mobility and changing requests of self-driving vehicles, and proposed a deep RL-based collaborative caching method to break the curse of high dimensionality in the state space of MDP, as well as decrease the redundant transmission delay.

The aforementioned studies are based on centralized methods that may not always be feasible in practice due to distributed and cooperative network topology. The fundamental problem is that these methods treat the other agents as a part of the environment and thus neglect the influence on the

TABLE I
NOTATIONS AND SYMBOLS

Notation	Explanation
\mathcal{N}	The set of all RSUs
\mathcal{F}	The index set of available contents
s_f	The size of content f
G_i	The limited available cache capacity of RSU i
\mathcal{U}	The set of all requested vehicles
ρ_f	The content popularity of content f
$T_{u,i}$	The content transmission delay between vehicle u and RSU i
$T_{i,j}$	The content transmission delay between RSU i and RSU j
$T_{i,N+1}$	The content transmission delay between RSU i and the MBS
$r_{u,i}$	The wireless downlink transmission rate between vehicle u and RSU i
$C_{u,i}$	The incurred transmission cost between vehicle u and RSU i
$C_{i,j}$	The incurred transmission cost between RSU i and RSU j
$C_{i,N+1}$	The incurred transmission cost between RSU i and the MBS
$\mathcal{L}_{f,i,j}$	The total overhead through different links
$a_{f,i,j}^t$	The caching state of local RSU i for the requested content f
$c_{f,i}^t$	The content replacement decision of RSU i in time slot t
$q_{i,\mathcal{F}}^t$	The content request arrived from the vehicles in time slot t

AQ:2

environment by other agents' behavior. Different from them, this paper investigates the edge caching strategy considering the content delivery and cache replacement in a distributed perspective. Specifically, the DMRE framwork is proposed to reduce the backhaul traffic in the system, where each agent can adaptively learn its best behavior in conjunction with other agents for intelligent caching. In addition, to further circumvent the issue by DMRE and make more expressive parametric models, another computationally efficient method (DeepDMRE) is also proposed, which approximates the Nash equilibria by constructing neural networks.

III. DISTRIBUTED VEHICULAR EDGE NETWORK ARCHITECTURE

This section presents the cooperative edge caching-supported IoVs architecture, including network infrastructure, content popularity pattern, as well as the process of content delivery and cache replacement. The main notations with their descriptions are listed in Table I.

A. Network Architecture

In this paper, a cooperative edge caching-supported IoVs architecture with an MBS and N small cells is considered, each of which is naturally equipped with an RSU, as shown in Fig. 1. The set of these RSUs is denoted by $\mathcal{N} = \{1, 2, \dots, N\}$. Considering the abundant cache capacity, it is assumed that MBS (denoted by $N + 1$ in the system) is connected to the Service Providers (SPs) and can cache all contents. Here, let $\mathcal{F} = \{1, 2, \dots, F\}$ represent the content library with total F available contents that are supported by the SPs, and s_f is the size of content f ($f \in \mathcal{F}$). Besides, each RSU i ($i \in \mathcal{N}$) is endowed with a limited available cache capacity of G_i , such that a small portion of popular contents can be cached in the RSUs to eliminate the redundant traffic.

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

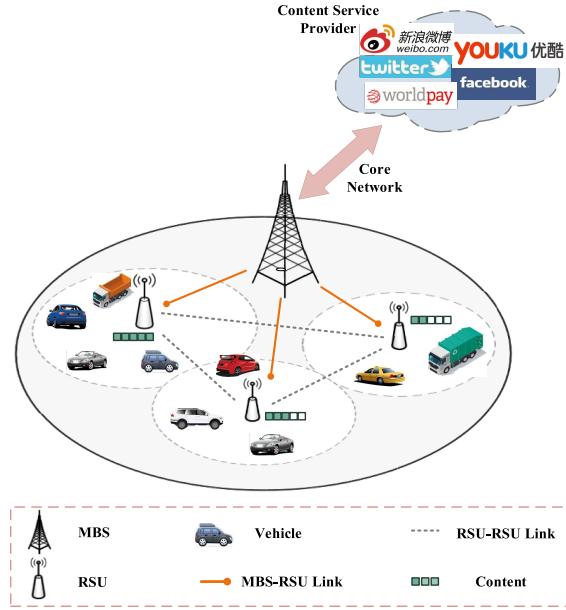


Fig. 1. Cooperative Edge Caching-supported IoVs Architecture.

All RSUs in a cluster form are connected to the remote MBS via wired lines. U vehicles (denoted as $\mathcal{U} = \{1, 2, \dots, U\}$) are randomly scattered in the wireless coverage scope and request various contents frequently via cellular links. Overlapping coverage between neighbouring cells is not considered, and these vehicles are assumed to be located in at least one small cell and will be served by the RSU therein; the set of vehicles in a small cell may change dynamically due to the vehicle mobility. Each RSU serves the local content requests within its coverage range. Meanwhile, the system is assumed to operate in a fixed length of time slots (i.e., time sequence) $t \in \{1, \dots, \Gamma\}$, where Γ denotes the finite time horizon. In each time slot, a vehicle can only request one content, while the caching decision of each RSU is also updated periodically.

As vehicles may exhibit different preferences for popular content, we introduce content popularity $\rho_f(f \in \mathcal{F})$ to reflect the content access requirements in the content library, i.e., ρ_f denotes the conditional probability of requesting content f among various vehicles' content requests. Inspired by existing works [27], [37], the content popularity is assumed to obey the Mandelbrot-Zipf distribution here. Then, the expected requesting probability for content f is obtained as:

$$\rho_f = \frac{(R_f + \tau)^{-\vartheta}}{\sum_{i \in \mathcal{F}} (R_i + \tau)^{-\vartheta}}, \quad \forall f \in \mathcal{F}, \quad (1)$$

where R_f is the rank of content f in the descending order of content popularity; $\vartheta > 0$ and τ are the skewness factor and plateau factor characterizing the distribution, respectively. Notably, due to the vehicle mobility, the local content popularity can be highly spatio-temporally varying over time slots and the requested vehicles in a cell, making caching strategy design more challenging.

Within the coverage of MBS, each RSU is able to cache various contents to meet the content access requirements from vehicles. Caching contents in RSUs enables the content

request to be accommodated in proximity, without duplicate downloading from remote MBS. Particularly, each RSU can determine where the content request is answered and which local cache should be replaced. Once a vehicle sends an access request within the range of the small cell, the local RSU will traverse its cache space firstly to check whether the desired content is cached in itself. If the content is not sought in its cache, the local RSU can either retrieve the content from the adjacent RSUs or download the content directly from the MBS, and deliver it to the requested vehicle later. Meanwhile, all RSUs may replace their cache with popular contents according to the content requests of each time slot. Underlying these is a sequential decision-making problem, which corresponds to making efficient content delivery and cache replacement decisions under constraints.

B. Content Delivery Model

$T_{u,i}$, $T_{i,j}$, and $T_{i,N+1}$ are used to denote the content transmission delay between vehicle u and RSU i , RSU i and RSU j , RSU i and MBS, respectively. It is worth noticing that the transmission delay of sending request identifier by the vehicle is neglected due to its tiny data size and the high link rate in most instances. Specifically, transmission delay of content f from any requested vehicle u ($u \in \mathcal{U}$) to its connected local RSU i can be calculated as $T_{u,i} = s_f / r_{u,i}$, where $r_{u,i}$ is the wireless downlink data rate between vehicle u and RSU i . Here, we assume that each RSU consists of multiple channels, and the allocated bandwidth to each channel is the same [38]. Therefore, considering the large-scale fading, the wireless downlink data rate $r_{u,i}$ is derived by

$$r_{u,i} = B_{u,i} \log_2 \left(1 + \frac{P_i g_{u,i}}{\sigma^2 + \sum_{v \in \mathcal{U} \setminus \{u\}} P_v g_{v,i}} \right), \quad (2)$$

where $B_{u,i}$ denotes the channel bandwidth, P_i denotes the transmission power of RSU i , σ^2 and $g_{u,i}$ are the white Gaussian noise power and the wireless propagation channel gain, respectively. Interference management is also considered for this wireless communication scenario. Furthermore, since the communications of MBS-RSU and RSU-RSU are via wired optical cables, the transmission delay $T_{i,N+1}$ and $T_{i,j}$ are set to constant values [35].

Delivering desired contents to vehicles will introduce extra transmission costs for caching nodes (cooperative RSUs or the MBS, denoted as $\mathcal{H} = \mathcal{N} \cup \{N+1\}$). As illustrated in Fig. 2, three components of transmission cost exist due to different cache hit situations. Here, $C_{u,i}$, $C_{i,j}$, and $C_{i,N+1}$ are used to denote the incurred transmission cost through RSU-Vehicle, RSU-RSU, and MBS-RSU links, respectively. Specifically, the incurred transmission cost through the RSU-Vehicle link is equal to the product of the link bandwidth, the content delivery delay, and the price of the RSU-Vehicle link's unit leased communication resource, which is expressed as:

$$C_{u,i} = B_{u,i} T_{u,i} \xi_{u,i}, \quad (3)$$

where $\xi_{u,i}$ represents the price of unit leased communication resource of wireless link. Similarly, the incurred transmission cost through RSU-RSU and MBS-RSU links can also be obtained by this way.

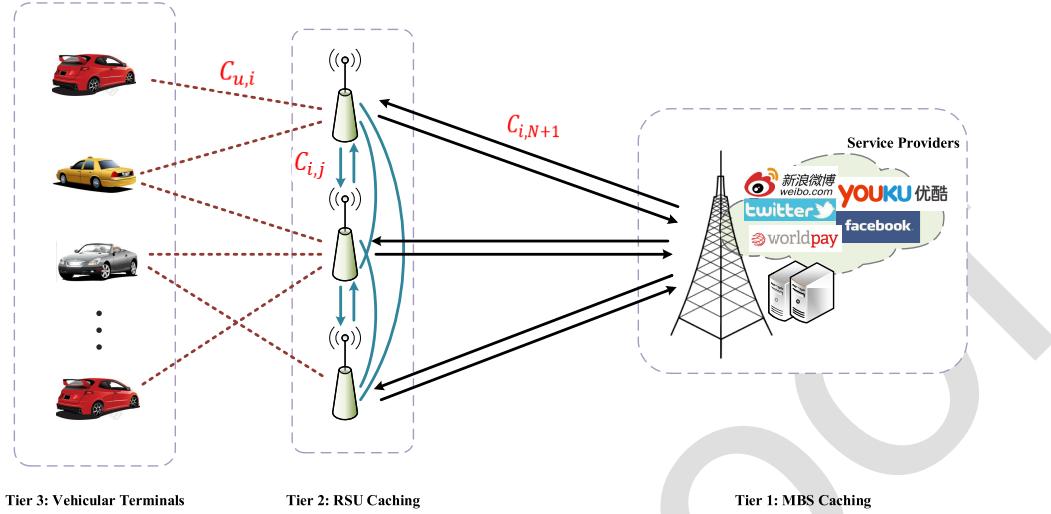


Fig. 2. Content delivery process and aggregate overhead in hierarchical vehicle edge networks.

Here, the binary decision variable $a_{f,i,j}^t \in \{0, 1\}$, $j \in \mathcal{H}$ is used to denote the caching state of any local RSU i for the requested content f in a given time slot t . When a content request occurs, the aggregate overhead of different links can be analyzed as follows:

- $a_{f,i,i}^t = 1$ indicates that content f is cached in the local RSU i at time slot t and can be delivered to the vehicle directly. In this case, the total overhead $\mathcal{L}_{f,i,i}$ is just the transmission cost $C_{u,i}$ from the local RSU to the vehicle;
- $a_{f,i,j}^t = 1$ ($j \in \mathcal{N} \setminus \{i\}$) indicates that content f is not cached in the local RSU i , but at least one of RSUs in the system has cached the desired content, thus the local RSU can inquire and obtain it from the cooperative RSU j . Then, the total overhead $\mathcal{L}_{f,i,j}$ consists of the transmission cost $C_{i,j}$ from the cooperative RSU j to the local RSU i and the transmission cost $C_{u,i}$;
- $a_{f,i,N+1}^t = 1$ indicates that there is no expected content in all cooperative RSUs, and the local RSU i has to forward the request identifier to MBS for processing at time slot t , i.e., the local RSU i downloads content f from MBS directly. In this way, the total overhead $\mathcal{L}_{f,i,N+1}$ consists of the transmission cost $C_{i,N+1}$ from MBS to the local RSU i and the transmission cost $C_{u,i}$.

C. Cache Replacement Model

In a diverse and dynamic edge caching scenario, it is indispensable to elaborate a wise content replacement strategy for efficiently managing the cache space. In order to meet the content requests that arrive later, each cooperative RSU should update the less popular contents in each operation phase to further improve the utilization and effectiveness of the whole cache.

Generally, the local RSU needs to determine whether to keep or replace its cache contents. When the newly requested contents arrive from other adjacent RSUs or MBS, some contents in the local RSU's existing cache may be evicted because of the finite cache space. As a result, the problem is whether and which contents should be replaced with the

new ones when the cache space is fully occupied. To this end, we represent the replacement control of RSU i by $c_{\mathcal{F},i}^t = \{c_{1,i}^t, \dots, c_{f,i}^t, \dots, c_{F,i}^t\}$ in time slot t , where $c_{f,i}^t \in \{0, 1\}$ ($f \in \mathcal{F}$) indicates whether and which content in RSU i should be replaced by the current content, e.g., $c_{f,i}^t = 1$ means that content f in RSU i needs to be replaced by the current arrived one, while $c_{f,i}^t = 0$ is the opposite. To efficiently store contents in RSUs, a utility-based replacement strategy will be introduced below.

IV. PROBLEM FORMULATION

This section formulates the optimization problem for IoV systems, and models the cache replacement process in an RSU as an MDP.

A. Objective Function

The RSUs in a cluster form can explore the potential of cooperation to fully utilize the cache resources. Considering the anticipated future utility, we intend to avoid myopic caching decisions while minimizing the long-term content access cost in the system. The corresponding problem is formulated as follows:

$$\begin{aligned} & \min_{a_{f,i,j}} \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \sum_{f=1}^F \sum_{i=1}^N \sum_{j=1}^{N+1} \rho_f a_{f,i,j}^t \mathcal{L}_{f,i,j}^t \\ & \text{s. t. } \quad C1 : a_{f,i,j}^t \in \{0, 1\}, \quad \forall f, \forall i, \forall j \\ & \quad C2 : a_{f,N+1,N+1}^t = 1, \quad \forall f \\ & \quad C3 : a_{f,i,j}^t \leq a_{f,j,j}^t, \quad \forall f, \forall i, \forall j \\ & \quad C4 : \sum_{j=1}^{N+1} a_{f,i,j}^t = 1, \quad \forall f, \forall i \\ & \quad C5 : \sum_{f=1}^F a_{f,i,i}^t s_f \leq G_i, \quad \forall i \end{aligned} \quad (4)$$

Here, $\lim_{T \rightarrow \infty} \frac{1}{T} (\dots)$ is the time averaged overhead of content delivery in the system. The meaning of the above constraints is

as follows: C1 guarantees the constraint of the binary caching decision's integer nature; C2 shows the MBS's sufficient cache capacity and guarantees that MBS has cached all available contents; C3 denotes that only the cooperative RSUs or MBS, which have cached the related content, can answer the content request; C4 ensures that the content requested by a particular vehicle can only be served by an RSU or MBS ultimately; C5 is used to promise that the cache usage of each RSU should be less than its cache capacity.

433 B. Problem and Challenges

434 In fact, it is challenging to solve the above problem
 435 in Eq. (4) directly since we have to obtain the optimal
 436 coordination of the caching decision variables $a_{f,i,j}^t$ ($f \in$
 $\mathcal{F}, j \in \mathcal{H} \setminus \{i\}$) in each time slot. However, the caching deci-
 437 sion variable $a_{f,i,j}^t$ of any RSU is binary and time-varying,
 438 which aggravates the difficulty by solving the problem at a
 439 time slot exhaustively. The system has to collect multitudinous
 440 network state information and make the global decision on
 441 cache replacement and content delivery for each RSU. More-
 442 over, we are devoted to a more practical case in which the prior
 443 information of the content request pattern is unknown. Thus,
 444 the optimization of the system is a combinatorial multi-armed
 445 bandit problem, and the objective function is undoubtedly
 446 NP-hard. Since the feasible set of the problem is not convex
 447 and the complexity is enormous, conventional methods using
 448 model-based heuristics or predefined rules may be unadaptable
 449 to make intelligent caching decisions under the dynamic
 450 system property.

452 C. Markov Decision Process

453 The MDP model is typically used to describe almost all
 454 sequential decision-making problems. Here, we model the
 455 optimization of cache replacement process in an RSU as an
 456 MDP. Three critical elements are identified as follows:

- **State Space:** The state in MDP is a space to reflect the IoV's environment. Accordingly, the state of an available RSU i is determined by the realization of the current caching situations and the request demand situations, which can be given as $z_t^i = \{a_{\mathcal{F},i,i}^t, q_{i,\mathcal{F}}^t\}$. As described earlier, the former $a_{\mathcal{F},i,i}^t = \{a_{1,i,i}^t, \dots, a_{f,i,i}^t, \dots, a_{F,i,i}^t\}$ denotes the current caching state respecting to the contents in RSU i . The latter $q_{i,\mathcal{F}}^t = \{q_{i,1}^t, \dots, q_{i,f}^t, \dots, q_{i,F}^t\}$ is the arrived content request state from vehicles in time slot t . Specifically, $q_{i,f}^t = 1$ means that at least one UE within local RSU i requests for content f in time slot t , $q_{i,f}^t = 1$ is the opposite. The above information will be aggregated as an input state in each time slot.
- **Action Space:** Once receiving a state in each time slot, the agent is responsible for deciding the caching and delivery strategy. Therefore, with the current state z_t^i , the action d_t^i involves two parts, $a_{\mathcal{F},i,j}^t$ and $c_{\mathcal{F},i}^t$, where matrix $a_{\mathcal{F},i,j}^t$ encodes the decision of RSU i for the requested contents, and the vector $c_{\mathcal{F},i}^t$ indicates the content replacement control in RSU i in slot t .

- **Reward Function:** Defining an appropriate reward function is vital since the reward value is the metric for each agent to evaluate the action quality. Generally speaking, the objective function is related to the immediate reward. Our objective in the related problem is to minimize the long-term content access cost in the system, which is the inverse goal of an agent that tries to achieve the maximum cumulative discounted reward. Thus, the reward function should be negatively correlated with the optimization objective. For this purpose, the immediate reward is defined as normalized $r(z_t^i, d_t^i) = e^{-\sum_f^F \sum_j^{N+1} \rho_f a_{f,i,j}^t L_{f,i,j}^t}$, where we utilize a negative exponential function to transform the problem legitimately. Furthermore, the reward value of an agent should satisfy the constraint conditions to ensure the validity of the results. Conversely, punishment negative will be incorporated into the reward if constraints are violated in Eq. (4).

V. MARL-EMPOWERED COOPERATIVE EDGE CACHING

Compared with dynamic programming methods, it is more suitable to adopt RL-based methods when the transition probability is ambiguous in MDP. This section extends MDP to a multi-agent system under the framework of a stochastic game, and proposes two MARL-based edge caching methods.

A. Markov Game Model

According to the MDP above, we have proposed a Q -learning based edge caching method in our previous work [7]. Q -learning is a model-free independent RL method based on value iteration. Each agent in Q -learning learns the settings separately through repeated interaction, and regards other agents as a part of the environment. Although single-agent properties are transferred directly to multi-agent contexts in Q -learning based method, the IoVs system is formulated as a multitude of subsystems where the cache replacement process in each available RSU is optimized individually. Nevertheless, as an independent caching node, each RSU should have its own caching strategy according to the corresponding content request and caching state, which may be very different between different RSUs. In addition, the caching strategies of agents are mutually influential in the distributed edge caching scenarios. Q -learning based method and their variants can only obtain the local sub-optimal solution, as they do not consider the influence on the environment by other agents when an agent interacts separately, i.e., each agent greedily makes decisions on its own benefit.

As far as these issues are concerned, we consider extending the MDP to a multi-agent system and formulating the cache replacement process as a Markov (a.k.a. Stochastic) Game (MG) model innovatively. In the MG model, due to the interaction between agents, the optimal decision of a single agent cannot ensure the optimal global solution of the system. Agents should not only observe their own immediate rewards and actions taken previously, but also those of others as well. Our objective is to obtain the best strategy for every agent. Therefore, it is necessary to make joint decisions among all

agents to avoid the deviation in a single decision, so as to learn the coordination and optimization of the whole system.

In MG model, agents execute actions at the same time, and the rewards of agents are generally arbitrarily related. Tuple $\{N, \mathcal{Z}, D_1, \dots, D_N, p, r_1, \dots, r_N, \gamma\}$ can characterize a standard formal definition of MG model, where N is the number of agents; \mathcal{Z} is the finite state space of the environment, each state $z_{t+1} \in \mathcal{Z}$ in the game consists of the individual states of all agents in the system; D_i stands for a discrete action space of agent i ($i = 1, \dots, N$); the joint action space of all agents is represented as $\mathcal{D} = D_1 \times \dots \times D_N$, and for notational convenience, we use $\vec{d}_t \doteq (d_t^1, d_t^2, \dots, d_t^N) \in \mathcal{D}$ to denote the current joint action of all agents in time slot t ; $p : \mathcal{Z} \times \mathcal{D} \times \mathcal{Z} \rightarrow [0, 1]$ is the state transition probability function, the joint action \vec{d}_t causes the transition from the game state z_t to the next new state z_{t+1} with a probability $p(z_{t+1}|z_t, \vec{d}_t)$, which should satisfy $\sum_{z_{t+1} \in \mathcal{Z}} p(z_{t+1}|z_t, \vec{d}_t) = 1$.

$\forall z_t \in \mathcal{Z}, \forall \vec{d}_t \in \mathcal{D}; r_i : \mathcal{Z} \times \mathcal{D} \rightarrow \mathbb{R}$ is the direct reward function of agent i , which evaluates the quality of state transition. Given the common state z_t , each agent i can perform an acceptable joint action $(d_t^1, \dots, d_t^i, \dots, d_t^N)$ and accumulate a reward $r_t^i(z_t, \vec{d}_t)$ immediately in new state z_{t+1} . Notably, the reward varies not only according to the current state and the action of agent i , but also the action choice of all other agents. Besides, state transitions obey the Markov property, where the sequence of events that precede the current state cannot determine the probability distribution of future states. An agent's reward and next state depend only upon the current state and the common joint decisions among all agents [39].

B. DMRE: Distributed MARL-Based Edge Caching Method

Since the distributed method is more effective in matching the edge caching system characteristics, in this part, Q -learning is expanded to the MG model, and DMRE is proposed to optimize the configuration of various caching strategies.

Indeed, the MG model can be seen as a set of matrix games associated with each state. In these games, each agent takes actions while conditioning on the behavior of others to satisfy the reward function, and $(\pi_1, \pi_2, \dots, \pi_N)$ is used to represent the joint strategy of N agents. Similar to Q -learning, the expected cumulative discounted reward of agent i is expressed by the state value function $V_i(z_t, \pi_1, \pi_2, \dots, \pi_N) = \mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t^i(z_t, \vec{d}_t)]$, where $\gamma \in (0, 1)$ is the discounting factor indicating the importance of the predicted future rewards. Each agent dedicates to learning a local π_i to maximize $V_i(z_t, \pi_1, \pi_2, \dots, \pi_N)$, meanwhile remaining robust to the actions of others over the remaining course of the game.

As the maximal state value $V_i(z_t, \pi_1^*, \dots, \pi_N^*)$ of an agent cannot be achieved by merely private strategy, and the return explicitly depends on the joint strategy of all agents, the concept of Nash equilibrium becomes crucial. Joint Nash equilibrium strategy $(\pi_1^*, \dots, \pi_i^*, \dots, \pi_N^*)$ can be considered

as an optimal solution of the multi-agent system, of which the game reaches the Nash equilibrium point and each agent's strategy π_i^* is the best response to the others. In other words, no agent can achieve a higher state value function by unilateral deviation, i.e., changing to any other strategy. Therefore, for $\forall z_t \in \mathcal{Z}, \forall i \in \mathcal{N}$, we formally have:

$$V_i(z_t, \pi_1^*, \dots, \pi_i^*, \dots, \pi_N^*) \geq V_i(z_t, \pi_1^*, \dots, \pi_i, \dots, \pi_N^*), \quad (5)$$

where $\forall \pi_i$ belongs to the set of strategies available to agent i .

Then, we desire a method to attain this equilibrium of the game without any prior knowledge about dynamics. Beforehand, we identify conditions that are more easily verifiable and extend the *Bellman Equation* to Nash equilibrium. Like Q -learning, we apply the dynamic programming principle to agent i 's reward $V_i(z_t, \pi_1^*, \dots, \pi_i^*, \dots, \pi_N^*)$ while the remaining policies $\pi_{\mathcal{N} \setminus \{i\}}^*$ are fixed, which will result in the temporal difference form as:

$$\begin{aligned} V_i(z_t, \pi_1^*, \dots, \pi_N^*) &= \max_{\pi_i} \left[r_t^i(z_t, \vec{d}_t) \right. \\ &\quad \left. + \gamma \sum_{z_{t+1} \in \mathcal{Z}} p(z_{t+1}|z_t, \vec{d}_t) V_i(z_{t+1}, \pi_i \cup \pi_{\mathcal{N} \setminus \{i\}}^*) \right]. \end{aligned} \quad (6)$$

Next, different from the single-agent Q -function² $Q_i(z_t, d_t)$, the Q -function in the multi-agent system varies to $Q_i(z_t, \vec{d}_t)$, which is interpreted to estimate the optimal expected sum of discounted rewards (state value function) for any individual agent i during the learning process, i.e., $V_i(z_t, \pi_1^*, \dots, \pi_N^*) = \max_{\pi_i} Q_i(z_t, \vec{d}_t)$. Specially, given the Nash equilibrium as a learning objective, $Q_i^*(z_t, \vec{d}_t)$ is referred to as a Nash Q -value for agent i when all agents follow a specified joint Nash equilibrium strategy from the next slot on. That is,

$$\begin{aligned} V_i(z_t, \pi_1^*, \dots, \pi_N^*) \\ = \text{Nash } Q_i(z_t, \vec{d}_t) \\ = \pi_1^*(z_t, d_t^1) \cdots \pi_N^*(z_t, d_t^N) \cdot Q_i^*(z_t, \vec{d}_t), \end{aligned} \quad (7)$$

where $\pi_1^*(z_t, d_t^1) \cdots \pi_N^*(z_t, d_t^N) \cdot Q_i^*(z_t, \vec{d}_t)$ is a scalar, and

$$\begin{aligned} Q_i^*(z_t, \vec{d}_t) \\ = r_t^i(z_t, \vec{d}_t) + \gamma \sum_{z_{t+1} \in \mathcal{Z}} p(z_{t+1}|z_t, \vec{d}_t) V_i(z_{t+1}, \pi_1^*, \dots, \pi_N^*). \end{aligned} \quad (8)$$

As a result, payoffs for each agent i are equal to Q -value at this point, and the Nash equilibrium is rewritten in the following form:

$$Q_i^*(z_t) \pi_1^*, \dots, \pi_i^*, \dots, \pi_N^* \geq Q_i^*(z_t) \pi_1^*, \dots, \pi_i, \dots, \pi_N^*. \quad (9)$$

²For ease of presentation, we have not distinguished the concept of “state-action function” and “ Q -function” in this paper.

At time slot t , each agent i executes its joint action under the current state. After that, it observes its immediate reward, all other agents' actions and rewards, as well as the new state z_{t+1} . As for agents' actions choice and Q -values update, the Nash equilibrium reward is adopted here to replace the maximum reward iteration in the single-agent method. Therefore, we rely on the stagewise approach to learn strategies separately for every stage game (i.e., a time slot game, defined by interim Q -values). While inequality (9) is satisfied, the stage game $(Q_t^1(z_{t+1}), \dots, Q_t^i(z_{t+1}), \dots, Q_t^N(z_{t+1}))$ will be formed with the Nash equilibrium strategies $\pi_1^*(z_{t+1}), \dots, \pi_N^*(z_{t+1})$ under certain restrictions to Q -values' domain. To achieve this, each learning agent has to keep track and maintain a model of other agents' Q -values, and then update its own Q -values during the interaction. The strategies that constitute a Nash equilibrium can be interpreted as the optimal behavior under the current system state, while different methods for selecting among multiple Nash equilibria will, in general, yield different updates. Existing studies [40], [41] show that there always exists at least one equilibrium point in stationary strategies, and the learning protocol provably converges if every stage game has saddle points or an optimum global; agents are mandated to update in terms of these. Ultimately, each agent i chooses its actions by calculating a Nash equilibrium, and updates the Q -value with the following iterative formula:

$$\begin{aligned} Q_{t+1}^i(z_t, \vec{d}_t) &= (1 - \beta_t) Q_t^i(z_t, \vec{d}_t) + \beta_t [r_t^i(z_t, \vec{d}_t) + \gamma \text{Nash } Q_t^i(z_{t+1})], \end{aligned} \quad (10)$$

where $\text{Nash } Q_t^i(z_{t+1})$ is the Nash equilibrium of agent i under new state, $\beta_t \in (0, 1)$ is the learning rate parameter, indicating how far the current estimate $Q_t^i(z_t, \vec{d}_t)$ is adjusted toward the update target $r_t^i(z_t, \vec{d}_t) + \gamma \text{Nash } Q_t^i(z_{t+1})$. Each Q -function can definitely converge to the Nash Q -value through repeated interaction when an appropriate β is designed. The agent can derive the Nash equilibrium and choose its actions accordingly.

C. Complexity Analysis and Parameter Approximation for DMRE

According to the description above, each agent i can gradually learn the equilibrium at each stage, which is determined by the joint actions of all agents in the system. That means the agent needs to learn N Q -values of all agents and derive strategies from them. Specifically, agent i will update $(Q^1, \dots, Q^i, \dots, Q^N)$ for all system states and joint actions, while these Q -values are updated and maintained internally by it. Hence, DMRE is often computationally expensive for equilibrium computation, and there may not be enough memory to maintain the corresponding values in large scale scenario. Consequently, it is necessary to analyze the computation complexity and further give simplification of the proposed method.

Let $|\mathcal{Z}|$ denote the number of system states, and $|D_i|$ be the size of agent i 's action space D_i . Assuming $|D_1| = \dots = |D_N| = |D|$, the total number of entries in Q^i is $|\mathcal{Z}| \cdot |D|^N$.

As the agent has to maintain N Q -tables, the total space requirement is $N|\mathcal{Z}| \cdot |D|^N$. Then, the proposed method is linear in the number of states, polynomial in the number of actions, but exponential in the number of agents in terms of space complexity. On the other hand, the training time of DMRE is dominated by the calculation of the Nash equilibrium. That is, each agent will search for an optimal solution according to the current solutions of other agents while the computation process is nonlinear if the number of agents is more than two [39]. Thus, due to the remarkable efficiency and generalization, parameter approximation is introduced to reduce the computation complexity for DMRE in this paper, which legitimately simplifies the training targets.

Specifically, elements can be controlled by sequences approximately in the space that is composed of Q -functions. Correspondingly, we introduce a set of parameters θ and approximate any agent i 's Q -function by using the updated parameter. The Q -function is re-expressed as:

$$\hat{Q}^i(z, \vec{d}, \theta) \approx \theta^T \phi_i = \sum_{j=1}^N \theta_j \phi_{ij}(z), \quad (11)$$

where $\theta = (\theta_1, \theta_2, \dots, \theta_n)^T$ and the given eigenfunction $\phi = (\phi_1(s), \phi_2(s), \dots, \phi_n(s))^T$.

Since $\hat{Q}^i(z, \vec{d}) \doteq \hat{Q}^i(z, \vec{d}, \theta)$ is possible when the error is smaller, we attempt to find a set of parameters θ for which the corresponding approximate values $(\hat{Q}_\theta^1, \dots, \hat{Q}_\theta^N)$ approximately satisfy Eq. (8). Thus, the loss function is interpreted as the distance between the approximate value and true value of Q -functions, which can be converted to:

$$\begin{aligned} \text{Loss}(\theta) &= \mathbb{E}_{z \sim p(z_{t+1}|z_t, \vec{d}_t)} \left[\left(Q^i(z, \vec{d}) - \hat{Q}^i(z, \vec{d}, \theta) \right)^2 \right] \\ &= \frac{1}{T} \sum_{t=1}^T \left[r_t^i + \gamma \text{Nash } \hat{Q}_t^i(z_{t+1}, \theta) - \hat{Q}_t^i(z_t, \vec{d}_t, \theta) \right]^2. \end{aligned} \quad (12)$$

According to the *Bellman Equation* of the state-action functions, Stochastic Gradient Descent (SGD) is performed to train θ towards the target value by minimizing the loss function, i.e. $\text{argmin}(\text{Loss}(\theta))$. Therefore, the optimal value is generated by an error as:

$$\theta_{t+1} \approx \theta_t - \frac{1}{2} \beta_t \nabla \left(r_t^i + \gamma \text{Nash } \hat{Q}_t^i(z_{t+1}, \theta) - \hat{Q}_t^i(z_t, \vec{d}_t, \theta) \right)^2. \quad (13)$$

We then take a semi-gradient method and ignore the derivative of parameter θ with respect to the unknown Nash Q value. Thus, the update formula of parameter θ in state-action function $\hat{Q}^i(z, \vec{d}, \theta)$ of agent i is written as follows:

$$\begin{aligned} \theta_{t+1} &= \theta_t + \beta_t \left[r_t^i + \gamma \text{Nash } \hat{Q}_t^i(z_{t+1}) \right. \\ &\quad \left. - \hat{Q}_t^i(z_t, \vec{d}_t, \theta_t) \right] \nabla \hat{Q}_t^i(z_t, \vec{d}_t, \theta_t). \end{aligned} \quad (14)$$

Instead of updating the Q -function, parameter approximation can make it more efficient by using the updated

Algorithm 1 Parameter Approximation for **DMRE**

Input: agent set \mathcal{N} , state space \mathcal{Z} , joint action space \mathcal{D} , learning rate β , discount factor γ , exploration factor ϵ .

```

1 Each agent  $i \in \mathcal{N}$  do
2   Initialize  $\forall z \in \mathcal{Z}, \forall d_i \in D_1, \hat{Q}^i(z, \vec{d}, \theta) \leftarrow 0$ , parameter  $\theta$ .
3   for each episode do
4     Set  $t = 0$ , obtain the initial state  $z_0$  by randomly caching.
5     for each slots of episode do
6       Derive an action  $d_t^i$  based on the  $\epsilon$ -greedy strategy in the current  $z_t$ .
7       for  $Q_i(z_t, \vec{d}_t) \neq Q_i^*(z_t, \vec{d}_t)$  do
8         Observe the rewards  $r_t^1, \dots, r_t^N$ , the joint actions  $\vec{d}_t$ , and the next state  $z_{t+1}$ ;
9         Update parameter  $\theta$  according to Eq. (14);
10        Compute  $Q^i(z_t, \vec{d}_t)$ :
11         $Q^i(z_t, \vec{d}_t, \theta) \leftarrow \theta^T \phi_i$ ;
12        Solve Nash equilibrium strategies under state  $z_{t+1}$ ;
13        Let  $t \leftarrow t + 1$ .
14     Execute cache strategy according to  $\pi_i^*$ .

```

parameter θ . Throughout the training process of seeking the Nash equilibrium point, the agent first updates the target θ by minimizing the corresponding loss function (12) via SGD. Then, θ can continually train the Q -function towards the Nash equilibrium point. All of these processes make the agent approach the Nash equilibrium point effectively. More details of the simplified DMRE³ are summarized in **Algorithm 1**, and relevant theoretical proofs of convergence and feasibility are presented at the end of this paper.

D. DeepDMRE: Distributed Deep MARL-Based Edge Caching Method

DMRE is enabled to represent and update the parameter by creating a look-up table, which is essentially a tabular-based method, even though parameter approximation is used to simplify the training target. Parameter update rule Eq. (14) in an agent relies on the repeated calculation of Nash $\hat{Q}_t^i(z_{t+1})$ over all states, which generally performs inefficiently in large-scale scenarios. To circumvent the issue and make more expressive parametric models, we incorporate the technical advantage of Deep- Q Network⁴ into DMRE, and further develop a computationally efficient method (DeepDMRE) with neural network-based Nash equilibria approximation.

³As the simplified DMRE has lower computation complexity and is easier to implement than the original version, we only use the “DMRE” to denote the simplified version in the following.

⁴Deep- Q Network is a classical deep RL method based on value iteration, in which neural networks are adopted to estimate the Q -function.

As mentioned above, we can introduce a set of parameters and approximate any agent’s Q -function using Eq. (11). In contrast, DeepDMRE defines a specific model for the function approximation. For each parameter θ , we decompose $\hat{Q}^i(z, \vec{d}, \theta)$ into two components:

$$\hat{Q}(z, \vec{d}, \theta) = \hat{V}(z, \theta) + \hat{A}(z, \vec{d}, \theta), \quad (15)$$

where $\hat{V}(z, \theta)$ and $\hat{A}(z, \vec{d}, \theta)$ are produced by various neural network components. $\hat{V}(z, \theta)$ is the state value function under state z , and $\hat{A}(z, \vec{d}, \theta)$ is the advantage function of executing action under state z . The advantage function represents the optimality gap between \hat{Q} and \hat{V} . In particular, each agent i ’s behaviour must be the best response to the others when the action execution constitutes a Nash equilibrium. At this point, the advantage function $\hat{A}(z, \vec{d}, \theta) = \text{Nash } Q_i(z_t, \vec{d}_t) - V_i(z_t, \pi_1^*, \dots, \pi_N^*) = 0$.

In contrast to DMRE, DeepDMRE maintains neural network models with experience replay buffers on each agent. The basic concept underlying DeepDMRE is centralized training and distributed inference. During training sample collection, the agent will store its observed experience tuple $y_t = (z_t, \vec{d}_t, z_{t+1}, r_t)$ into the replay buffer, which involves its current state and available action set. The arrived experience samples are used to train the parameters of the function approximation in neural networks via SGD. Practically, this replaying enables the agent to randomly extract a minibatch of previous experience samples from the replay buffer for learning at each iteration. Empirical studies [15], [20], [33] have proved that the experience replay mechanism can enhance sample efficiency and break the correlation among data training.

Moreover, we generalize and adopt the model hypothesis proposed in [40], which considers a second order Taylor expansion $\hat{Q}^i(z, \vec{d}, \theta)$ in action around the Nash equilibrium. Together with the hypothesis, the advantage function for each agent i can be approximately written in a linear-quadratic form as:

$$\begin{aligned} \hat{A}_i(z, \vec{d}, \theta) &= \left(\vec{d}_{\mathcal{N} \setminus \{i\}} - \bar{\mu}_{\mathcal{N} \setminus \{i\}}^\theta \right)^\top \Psi_i^\theta(z) \\ &\quad - \left(\frac{d_i - \mu_i^\theta}{\vec{d}_{\mathcal{N} \setminus \{i\}} - \bar{\mu}_{\mathcal{N} \setminus \{i\}}^\theta} \right)^\top P_i^\theta(z) \left(\frac{d_i - \mu_i^\theta}{\vec{d}_{\mathcal{N} \setminus \{i\}} - \bar{\mu}_{\mathcal{N} \setminus \{i\}}^\theta} \right) \end{aligned} \quad (16)$$

where the block matrix $\Psi_i^\theta(x)$ and $P_i^\theta(x)$ are all deterministic probability distributions. Interested readers can find detailed information in [40].

This model hypothesis implicitly suggests that each agent’s Q -function can be approximately expressed as a linear-quadratic function of the actions. Each $\hat{Q}^i(z, \vec{d}, \theta)$ is a concave function of d_i , guaranteeing that Nash $\hat{Q}_t^i(z_{t+1})$ is bijective. Then, we can model $\hat{V}(z, \theta)$, $\bar{\mu}^\theta$, and relevant probability distributions via constructed neural networks, rather than modelling $\hat{Q}^i(z, \vec{d}, \theta)$. The Nash equilibrium is achieved when $\vec{d}^* = \bar{\mu}^\theta$, and the advantage function will become 0 simultaneously. Similar to the derivations in Eq. (7), we can simplify expressions of the value function and equilibrium

801 strategy as:

$$\hat{V}(z, \theta) = \text{Nash } \hat{Q}^i(z, \vec{d}, \theta) \quad (17)$$

802 and

$$\vec{\mu}(z) = \text{argmax } \text{Nash } \hat{Q}^i(z, \vec{d}, \theta) \quad (18)$$

803 Consequently, the loss function in Eq. (12) becomes more
804 tractable through this simplification. For a minibatch of J
805 sample, it can be converted as:

$$\begin{aligned} 806 \mathcal{L}(\theta) &= \frac{1}{J} \sum_{j=1}^J \mathcal{L}_j(\theta) \\ 807 &= \frac{1}{J} \sum_{j=1}^J \left(r_t^i + \gamma \hat{V}_t^i(z_{t+1}, \theta) - \hat{V}(z, \theta) - \hat{A}(z, \vec{d}, \theta) \right)^2. \end{aligned} \quad (19)$$

808 where each sample observation $y_j = (z_j, \vec{d}_j, z_{j+1}, r_j)$ is used
809 to train the parameter θ by minimizing the $\mathcal{L}_j(\theta)$.

810 Furthermore, instead of applying SGD with back-
811 propagation, we employ an actor-critic-based variant on the
812 loss function, enabling more stable and efficient convergence
813 processes. Specifically, there are N actor networks and one
814 centralized critic network, where each RSU adopts one of the
815 actor networks to execute its inference, and the MBS acts as a
816 centralized critic network to train the model and evaluate the
817 overall caching state. Since the value function can be modeled
818 independently through the decomposition in Eq. (15), an actor-
819 network is performed for parameters update by minimizing
820 the loss function in Eq. (19) over parameters θ . Notably, the
821 parameter θ gets separated as $\theta = (\theta_V, \theta_A)$ in this update rule,
822 where sub-parameter θ_V is used to model $\hat{V}(z, \theta_V)$, and sub-
823 parameter θ_A is used to model $\hat{A}(z, \vec{d}, \theta_A)$. Accordingly, the
824 training objective is to update these parameters by minimizing
825 the redefined loss function:

$$\mathcal{L}(\theta) = \frac{1}{J} \sum_{j=1}^J \mathcal{L}_j(\theta) = \frac{1}{J} \sum_{j=1}^J \hat{\mathcal{L}}(y_j, \theta_V, \theta_A), \quad (20)$$

826 where each agent minimizes $\mathcal{L}(\theta)$ by alternating between
827 minimization in sub-parameters θ_V and θ_A . Here, the loss
828 function of an individual sample $y_j = (z_j, \vec{d}_j, z_{j+1}, r_j)$
829 corresponds to the error in Eq. (8), which is expressed as:

$$\begin{aligned} 830 \hat{\mathcal{L}}(y_j, \theta_V, \theta_A) \\ 831 &= \left(r(z_j, \vec{d}_j) + \gamma \hat{V}(z_{j+1}, \theta_V) - \hat{V}(z_j, \theta_V) - \hat{A}(z_j, \vec{d}_j, \theta_A) \right)^2. \end{aligned} \quad (21)$$

832 More details of DeepDMRE are summarized in
833 **Algorithm 2**. With the locally linear-quadratic form of
834 the advantage function, we minimize the loss function in
835 Eq. (21) over sub-parameters θ_V and θ_A by actor-critic-based
836 iterative optimization and batch sampling. On the other hand,
837 we also apply and modify the ϵ -greedy based action
838 selection strategy for ensuring adequate exploration, where
839 ϵ is a decreasing parameter to achieve a tradeoff between
840 exploitation and exploration.

Algorithm 2 DeepDMRE: Distributed Deep MARL Based Edge Caching

841 **Input:** agent set \mathcal{N} , state space \mathcal{Z} , joint action space
842 \mathcal{D} , learning rate β , discount factor γ ,
843 exploration factor ϵ , experience replay buffer
844 \mathcal{M}_i , minibatch size J .

845 1 **Each agent** $i \in \mathcal{N}$ **do**

846 2 **Initialize** $\forall z \in \mathcal{Z}, \forall d_i \in D_i, \hat{Q}^i(z, \vec{d}, \theta) \leftarrow 0$, replay
847 buffer \mathcal{M}_i , neural network parameter (θ_V, θ_A) .

848 3 **for** each episode **do**

849 4 Set $t = 0$, obtain the initial state z_0 by randomly
850 caching.

851 5 **for** each slots of episode **do**

852 6 Derive an action $\vec{d} \leftarrow \vec{\mu}^\theta$ based on the ϵ -greedy
853 strategy in the current z_t .

854 7 **for** $Q_i(z_t, \vec{d}) \neq Q_i^*(z_t, \vec{d})$ **do**

855 8 Observe the rewards r_t^1, \dots, r_t^N , the
856 actions \vec{d}_t , and the next state z_{t+1} ;

857 9 Store observed experience tuple
858 $y_t = (z_t, \vec{d}_t, z_{t+1}, r_t)$ into the replay
859 buffer.

860 10 **if** $t \% \text{updateFrequency} == 0$ **then**

861 11 Randomly extract a minibatch of J
862 experiences $\{y_j\}_{j=1}^J$ from replay
863 buffer \mathcal{M}_i .

864 12 Update sub-parameters θ_V and θ_A by
865 minimizing $\mathcal{L}(\theta)$ as Eq. (20);

866 13 Perform a an actor-critic-based variant on
867 the loss function with respect to the
868 sub-parameters θ by $\frac{\partial \text{Loss}(\theta)}{\partial \theta}$;

869 14 Let $t \leftarrow t + 1$.

870 15 Update (θ_V, θ_A) .

VI. PERFORMANCE EVALUATION

871 This section conducts extensive simulations to assess the
872 performance of the proposed methods. Specifically, the su-
873 periority of the proposed methods is demonstrated over both
874 rule-based and learning-based benchmarks.

A. Simulation Setup

875 This paper considers an IoVs architecture composed of
876 one MBS and 5 RSUs, in which the coverage area of each
877 RSU is 200 m. RSUs are located at the center of each
878 region to serve the corresponding content requests, and the
879 initial cache capability of each RSU is limited to 100 MB.
880 Notably, the model can be easily extended to the case of
881 inconsistent RSU cache capacity. Besides, 10, 000 contents
882 are generated, and the size of each content is within the
883 range of [0.5, 1.5] MB. The content popularity follows an
884 MZipf distribution with $\theta = 0.73$. According to [4] and [35],
885 we set the channel gain as $30.6 + 36.7 \log_{10}(d)$, and the
886 noise power σ^2 as -95 dBm. The transmission delays between
887 MBS and RSU, and cooperative RSUs are 80 ms and 10 ms,

TABLE II
PARAMETERS SETUP

Parameter	Definition	Value
N	Default RSU number	5
T_2	Delay between any cooperative RSUs	10ms
T_3	Delay between the MBS and an RSU	80ms
P_i	Transmission power of RSUs	38dBm
B	Bandwidth of wireless links	10MHz
W	Bandwidth of optical fibre links	20MHz
σ^2	Power of background noise	-95 dBm
s_f	Size of requested contents	[0.5, 1.5] MB

865 respectively. Here, the price of the unit leased communication
866 resources (\$ / MB) is set as 0.005 for wireless links and
867 0.009 for optical links, while their bandwidths are 10 MHz
868 and 20 MHz, respectively. Meanwhile, we set the capacity of
869 experience replay buffer $\mathcal{M} = 500$, minibatch size $J = 32$,
870 and initial exploration factor $\epsilon = 0.05$. The learning rate
871 parameters for the actor-network and the critic-network are
872 set as 10^{-3} and 3×10^{-4} , respectively. All experiments are
873 run on 24 core - 2.4GHz Intel Xeon E5-2650 processor and
874 256GB RAM. The main parameters are listed in Table II.

875 For performance comparison, the following four benchmark
876 caching methods are introduced:

- 877 1) Least Frequently Used (LFU): The content with the least
 878 requested times will be replaced firstly when an RSU's
 879 cache capacity is full.
- 880 2) Least Recently Used (LRU): The content with the longest unused time will be replaced firstly when an
 881 RSU's cache capacity is full.
- 882 3) Q-learning [7]: An independent reinforcement learning
 883 method, the process of cache replacement in each available
 884 RSU is optimized individually.
- 885 4) Informed Upper Bound (IUB) [13]: An omniscient
 886 method which assumes all of the requesting information
 887 is perfectly known in advance. It provides an upper
 888 bound of other methods and is usually hard to achieve
 889 in practice.

890 Besides, the following metrics are used to evaluate the
891 performance of different methods quantitatively. 1) Content
892 access cost: the average total access cost of all requests in
893 each time slot. 2) Edge hit rate: the sum of requests served
894 by RSUs divided by the total number of requests. 3) Average
895 delay: the average access time for all requests in each time
896 slot.

898 B. Performance Comparison

899 1) *The convergence performance*: We plot the convergence
900 curve for three RL-based methods in Fig. 3, where the
901 number of contents and the cache capability of RSU are
902 fixed at 10,000 and 100MB, respectively. We can observe
903 that for these methods, each episode's total content access
904 cost decreases rapidly with the increase of the number of
905 episodes. The DeepDMRE converges to a relatively stable
906 value after running about 1800 episodes. Overall, it performs
907 best as it achieves a faster running time and obtains the same
908 convergence value as the original DMRE. On the other hand,

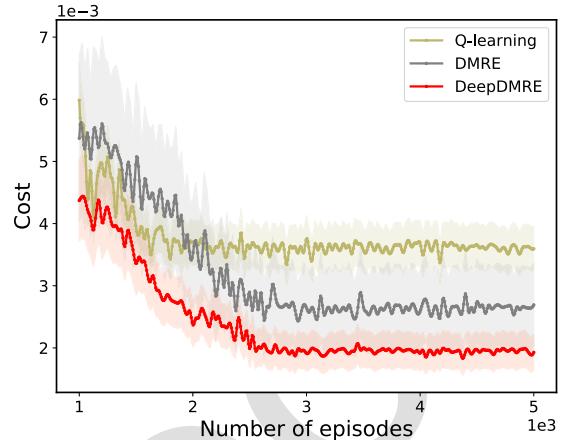


Fig. 3. Content access cost versus the number of episodes.

909 although the convergence speed of distributed MARL-based
910 methods (DMRE and DeepDMRE) are slightly slower than
911 that of the independent Q-learning based method, they can
912 reduce the content access cost up to 50% during the training
913 episodes.

914 2) *The Impact of Cache Capability of RSUs*: Then,
915 we explore the impact of different cache capabilities of RSUs
916 on the performance of different methods, as shown in Fig. 4.
917 The cache capacity of each RSU is changed from 100 MB
918 to 1000 MB and the number of contents is 10,000 in this case.
919 As expected, IUB performs best. However, it is impractical
920 since future information cannot be perfectly known. Among
921 the other caching methods, we observe that as the cache
922 capacity of RSUs increases, DeepDMRE consistently achieves
923 better performance with respect to the content access cost and
924 average delay. Specifically, DeepDMRE achieves the lowest
925 average delay of 33ms in the initial stage, and the content
926 access cost is reduced by about 47% and 26% compared with
927 independent Q-learning and DMRE, without mentioning the
928 simple rule-based methods LFU and LRU. In addition, for
929 all caching methods, the content access cost and the average
930 delay decrease with the increase of the cache capacity, which
931 is reasonable. The larger cache capacity enables the RSU to
932 cache more contents at the same time so that RSUs can meet
933 most content requests.

934 Similarly, the increase of cache capacity also has a positive
935 impact on edge hit rate. From Fig. 4(b), we can find that
936 the edge hit rate exhibits an increasing trend for all caching
937 methods when we increase the cache capacity of each RSU,
938 especially for DMRE and DeepDMRE, which are very close to
939 the value of IUB. This is because they considers the influence
940 of environments by other agents, and takes joint decisions to
941 learn the coordination and optimization of the whole system.
942 Notably, DMRE and DeepDMRE will generally sacrifice some
943 local hit rate and share cache capacity to serve requests from
944 neighboring RSUs. This tradeoff benefits much by reducing
945 the duplicate content transmission since MBS fetching is
946 largely avoided. Especially, DeepDMRE outperforms DMRE,
947 Q-learning, LFU, and LRU. For example, the edge hit rate is
948 improved by roughly 5%, 19%, 40%, and 35%, respectively,
949 when the cache capacity reaches 1,000 MB. At this moment,
950 content replacement processes may rarely occur in the RSUs.

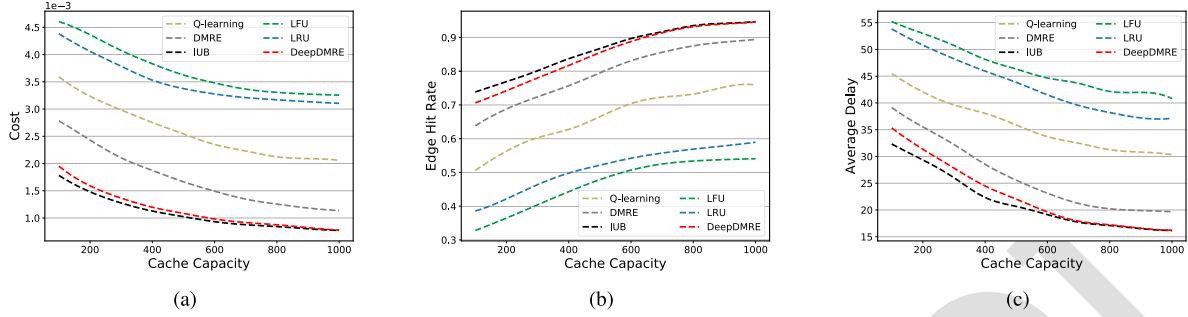


Fig. 4. (a) Content access cost versus the cache capacity of RSUs. (b) Edge hit rate versus the cache capacity of RSUs. (c) Average Delay versus the cache capacity of RSUs.

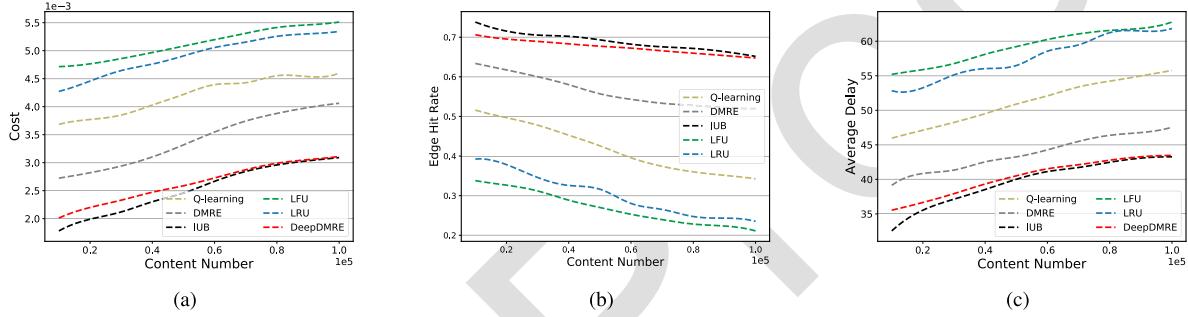


Fig. 5. (a) Content access cost versus the number of contents. (b) Edge hit rate versus the number of contents. (c) Average Delay versus the number of contents.

951 3) *The Impact of the Number of Contents in the System:*
952 This part changes the number of contents to evaluate the
953 performance of different methods. The number of contents is
954 varied from 10,000 to 100,000 and the initial cache capacity
955 of RSUs is set as 100 MB. As illustrated in Fig. 5, the content
956 access cost and the average delay of all methods increase
957 slightly as the number of contents increases. The increasing
958 trend does indicate that when the cache capacity is finite,
959 more contents will naturally cause frequent cache replacement,
960 as more contents need to be cached in RSUs now. As expected,
961 DeepDMRE performs better with a different numbers of
962 contents, except IUB. From Fig. 5(a) and (c), we can find
963 that even when there is a massive number of contents (i.e.,
964 100,000), DeepDMRE can still reduce the content access cost
965 by 11%, 16%, 28%, 34%, and the average delay by 15%, 21%,
966 33%, 37% compared to DMRE, Q-learning, LRU, and LFU,
967 respectively, which infers the effectiveness of DeepDMRE
968 under different numbers of contents.

969 Besides, the edge hit ratio of all methods decreases with
970 the increase of the number of contents. Among them, edge
971 hit rate of DeepDMRE is 64% when the number of contents
972 reaches 100,000. In contrast, LFU and LRU still perform
973 worst for this moment, and the edge hit rate is only 21%
974 and 23%, respectively. This situation may occur because LFU
975 and LRU learn only from one-step past and operate based
976 on simple rules, while RL-based edge caching methods can
977 be derived from the observed historical content demands and
978 concentrate more on the reward that agents can earn rather
979 than users' requests. Furthermore, the variation of edge hit
980 rate also confirms the analysis above. That is, DeepDMRE
981 improves the cache resources utilization significantly.

982 To summarize, quantitative results validate the effectiveness
983 of DMRE and DeepDMRE under different scenarios.
984 As observed, they not only reduce the content access cost,
985 but also achieve a desirable edge hit rate and average
986 delay.

VII. CONCLUSION

987 This paper first considered a cooperative edge caching sup-
988 ported IoVs architecture, which uses the cooperative caching
989 between multiple RSUs and MBS to reduce the duplicate
990 content transmission in the system. Then, the MDP was
991 extended to the context of a multi-agent system, and the
992 corresponding combinatorial multi-armed bandit problem was
993 further tackled based on the framework of stochastic games.
994 Specifically, a **Distributed MARL-based Edge** caching method
995 (DMRE) was proposed firstly, where each agent can adaptively
996 learn its best behavior in conjunction with other agents for
997 intelligent caching. Meanwhile, we attempted to reduce the
998 computation complexity of DMRE by parameter approxima-
999 tion, which legitimately simplifies the training targets. How-
1000 ever, DMRE is enabled to represent and update the parameter
1001 by creating a lookup table, essentially a tabular-based method,
1002 which generally performs inefficiency in large-scale scenarios.
1003 To circumvent the issue and make more expressive parametric
1004 models, we combined the technical advantage of the Deep-Q
1005 Network into DMRE, and further developed a computa-
1006 tionally efficient method, named DeepDMRE, by construct-
1007 ing neural networks for approximating the Nash equilibria.
1008 Simulation results demonstrated that DeepDMRE signifi-
1009 cantly outperforms other benchmark methods in different
1010 scenarios.

APPENDIX

This part gives the relevant theoretical proofs of convergence and the feasibility of the simplified DMRE as follows.

A. Convergence Analysis

Since the original DMRE has a provably convergence performance with restrictive conditions [39], [41], the convergence of the simplified DMRE can accordingly be proved by demonstrating the effectiveness of SGD. Therefore, as θ (which minimizes the corresponding loss function) is gradually obtained via SGD, the state-action function (in the form of $\hat{Q}^i(z, \vec{d}, \theta)$) will be updated towards the Nash equilibrium point.

For any state-action function $Q^k (k \in \mathcal{N})$ maintained internally by a learning agent, we can define $f_k(\theta) = (Q^k(z, \vec{d}) - \hat{Q}^k(z, \vec{d}, \theta))^2$. Then, the corresponding loss function is written as:

$$L(\theta) \triangleq Loss(\theta) = \frac{1}{N} \sum_{k=1}^N f_k(\theta). \quad (22)$$

As described in Section V-C, we train the parameter θ towards the target value by minimizing the loss function at each iteration. Then, owing to the numerical convexity of the loss function, its gradient $\nabla L(\theta)$ satisfies the *Lipschitz Continuity* with respect to the constant G , so that for all θ_t and θ_{t+1} . We have:

$$\|\nabla L(\theta_t) - \nabla L(\theta_{t+1})\| \leq G \|\theta_t - \theta_{t+1}\|. \quad (23)$$

For a series of progressively decreasing learning rates β_t , the iterative formula of parameter θ can be expressed as:

$$\theta_{t+1} = \theta_t - \beta_t \nabla L(\theta_t). \quad (24)$$

Meanwhile, each θ in the training process should satisfy $\max_k \|\nabla f_k(\theta_t)\| \leq B \|\nabla L(\theta_t)\|$, where B is a constant. That is to say, the optimal θ for loss function is supposed to be a stagnation point at the same time. We denote the error during the iteration process as e_t , then the above Eq. (24) can be rewritten with the following form:

$$\theta_{t+1} = \theta_t - \beta_t (\nabla L(\theta_t) + e_t), \quad (25)$$

where the error $e_t = \nabla f_k(\theta_t) - \nabla L(\theta_t)$.

Additionally, since the mean of the error is equal to 0, we therefore get:

$$E[e_t] = E[\nabla f_k(\theta_t) - \nabla L(\theta_t)] = E[\nabla f_k(\theta_t)] - \nabla L(\theta_t) = 0. \quad (26)$$

Furthermore, based on the derivations above, we can easily get:

$$\begin{aligned} & E[\|e_t\|^2] \\ &= E[\|\nabla f_k(\theta_t) - \nabla L(\theta_t)\|^2] \\ &= E[\|\nabla f_k(\theta_t)\|^2 - 2 \langle \nabla f_k(\theta_t), \nabla L(\theta_t) \rangle + \|\nabla L(\theta_t)\|^2] \\ &= E[\|\nabla f_k(\theta_t)\|^2 - 2 \langle E[\nabla f_k(\theta_t)], \nabla L(\theta_t) \rangle + \|\nabla L(\theta_t)\|^2] \end{aligned}$$

$$\begin{aligned} &= \frac{1}{N} \sum_{k=1}^N [\|\nabla f_k(\theta_t)\|^2 - \|\nabla L(\theta_t)\|^2] \\ &\leq (B^2 - 1) \|\nabla L(\theta_t)\|^2. \end{aligned} \quad (27)$$

Then, for inexact gradients, inequality (23) can be rewritten by incorporating Eq. (24) with it:

$$L(\theta_{t+1}) \leq L(\theta_t) + \langle \nabla L(\theta_t), \theta_{t+1} - \theta_t \rangle + \frac{G}{2} \|\theta_{t+1} - \theta_t\|^2. \quad (28)$$

Also, we substitute error $e_t = \nabla f_k(\theta_t) - \nabla L(\theta_t)$ and Eq. (24) into inequality (28), and further obtain:

$$\begin{aligned} L(\theta_{t+1}) &\leq L(\theta_t) - \beta_t \langle \nabla L(\theta_t), \nabla L(\theta_t) + e_t \rangle \\ &\quad + \frac{\beta_t^2 G}{2} \|\nabla L(\theta_t) + e_t\|^2 \\ &= L(\theta_t) - \beta_t \left(1 - \frac{\beta_t G}{2}\right) \|\nabla L(\theta_t)\|^2 \\ &\quad - \beta_t (1 - \beta_t G) \langle \nabla L(\theta_t), e_t \rangle + \frac{\beta_t^2 G}{2} \|e_t\|^2. \end{aligned} \quad (29)$$

As the learning rate β_t tends to be infinitesimal, we take the expectation of e_t on both sides of inequality (29). Then, a new inequality (30), shown at the top of the next page, is derived based on (26) and (27), as is shown at the bottom of this paper. It shows that SGD can obtain the expected parameter θ_t as the learning rate is small enough (satisfy $0 < \beta_t < \frac{2}{GB^2}$). At this point, we can get the corresponding state-action function by the obtained parameter θ_t , the effectiveness of SGD is therefore demonstrated. According to the previous description, the convergence of the simplified DMRE is proved theoretically.

B. Feasibility Analysis

To reduce the computation complexity, parameter approximation is introduced for DMRE in this paper, which legitimately simplifies the training targets. However, a potential error m may exist in the simplified DMRE, which is written as:

$$m = Q^i(z, \vec{d}) - \hat{Q}^i(z, \vec{d}, \theta) = Q^i(z, \vec{d}) - \sum_{j=1}^n \theta_j \phi_{ij}(z). \quad (31)$$

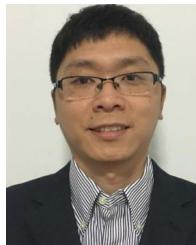
According to inequality (30), the target parameter which minimizes the loss function will be obtained via SGD throughout the training process. Thus, the error possibly existing in the simplified DMRE is contingent on the performance of SGD actually. That is to say, the optimal parameter value can minimize the corresponding error. On the other hand, an agent needs to learn N Q -values of all agents and derive strategies from them in the original DMRE, while these Q -values are updated and maintained internally by it. This pattern makes it obvious to get trouble in the curse of dimensionality with the increase of agents. Moreover, the computation complexity increases considerably, as well as the running time. To avoid these bottlenecks, we further give an simplification of DMRE, which legitimately transforms the training targets through parameter approximation in Eq. (14).

$$\begin{aligned}
E[L(\theta_{t+1})] &\leq L(\theta_t) - \beta_t \left(1 - \frac{\beta_t G}{2}\right) \|\nabla L(\theta_t)\|^2 - \beta_t (1 - \beta_t G) \langle \nabla L(\theta_t), E[e_t] \rangle + \frac{\beta_t^2 G (B^2 - 1)}{2} \|\nabla L(\theta_t)\|^2 \\
&\leq L(\theta_t) - \beta_t \left(1 - \frac{\beta_t G}{2}\right) \|\nabla L(\theta_t)\|^2 + \frac{\beta_t^2 G (B^2 - 1)}{2} \|\nabla L(\theta_t)\|^2 \\
&\leq L(\theta_t) - \beta_t \left(1 - \frac{\beta_t G B^2}{2}\right) \|\nabla L(\theta_t)\|^2.
\end{aligned} \tag{30}$$

REFERENCES

- [1] P. Arthurs, L. Gillam, P. Krause, N. Wang, K. Halder, and A. Mouzakitis, "A taxonomy and survey of edge cloud computing for intelligent transportation systems and connected vehicles," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 7, pp. 6206–6221, Jul. 2022.
- [2] K. Jiang, C. Sun, H. Zhou, X. Li, M. Dong, and V. C. M. Leung, "Intelligence-empowered mobile edge computing: Framework, issues, implementation, and outlook," *IEEE Netw.*, vol. 35, no. 5, pp. 74–82, Sep. 2021.
- [3] Q. Luo, S. Hu, C. Li, G. Li, and W. Shi, "Resource scheduling in edge computing: A survey," *IEEE Commun. Surveys Tuts.*, vol. 23, no. 4, pp. 2131–2165, 4th Quart., 2021.
- [4] J. Zhang and K. B. Letaief, "Mobile edge intelligence and computing for the Internet of Vehicles," *Proc. IEEE*, vol. 108, no. 2, pp. 246–261, Feb. 2019.
- [5] H. Zhou, M. Li, N. Wang, G. Min, and J. Wu, "Accelerating deep learning inference via model parallelism and partial computation offloading," *IEEE Trans. Parallel Distrib. Syst.*, vol. 34, no. 2, pp. 475–488, Feb. 2023.
- [6] Y. Zhang, X. Lan, J. Ren, and L. Cai, "Efficient computing resource sharing for mobile edge-cloud computing networks," *IEEE/ACM Trans. Netw.*, vol. 28, no. 3, pp. 1227–1240, Jun. 2020.
- [7] K. Jiang, H. Zhou, D. Zeng, and J. Wu, "Multi-agent reinforcement learning for cooperative edge caching in Internet of Vehicles," in *Proc. IEEE 17th Int. Conf. Mobile Ad Hoc Sensor Syst. (MASS)*, Dec. 2020, pp. 455–463.
- [8] J. Cui, Y. Liu, Z. Ding, P. Fan, A. Nallanathan, and L. Hanzo, "Next-generation mm-Wave small-cell networks: Multiple access, caching, and resource management," *IEEE Veh. Technol. Mag.*, vol. 15, no. 1, pp. 46–53, Mar. 2020.
- [9] L. Su and V. K. N. Lau, "Data and channel-adaptive sensor scheduling for federated edge learning via over-the-air gradient aggregation," *IEEE Internet Things J.*, vol. 9, no. 3, pp. 1640–1654, Feb. 2022.
- [10] H. Li, K. Ota, and M. Dong, "Deep reinforcement scheduling for mobile crowdsensing in fog computing," *ACM Trans. Internet Technol.*, vol. 19, no. 2, pp. 1–18, May 2019.
- [11] H. Zhou, Z. Zhang, D. Li, and Z. Su, "Joint optimization of computing offloading and service caching in edge computing-based smart grid," *IEEE Trans. Cloud Comput.*, early access, Apr. 12, 2022, doi: 10.1109/TCC.2022.3163750.
- [12] W. Zhuang, Q. Ye, F. Lyu, N. Cheng, and J. Ren, "SDN/NFV-empowered future IoV with enhanced communication, computing, and caching," *Proc. IEEE*, vol. 108, no. 2, pp. 274–291, Feb. 2020.
- [13] W. Jiang, G. Feng, S. Qin, T. S. P. Yum, and G. Cao, "Multi-agent reinforcement learning for efficient content caching in mobile D2D networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 3, pp. 1610–1622, Mar. 2019.
- [14] H. Zhou, T. Wu, H. Zhang, and J. Wu, "Incentive-driven deep reinforcement learning for content caching and D2D offloading," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 8, pp. 2445–2460, Aug. 2021.
- [15] H. Zhang, M. Huang, H. Zhou, X. Wang, N. Wang, and K. Long, "Capacity maximization in RIS-UAV networks: A DDQN-based trajectory and phase shift optimization approach," *IEEE Trans. Wireless Commun.*, vol. 22, no. 4, pp. 2583–2591, Apr. 2023, doi: 10.1109/TWC.2022.3212830.
- [16] H. Zhang, N. Yang, W. Huangfu, K. Long, and V. C. M. Leung, "Power control based on deep reinforcement learning for spectrum sharing," *IEEE Trans. Wireless Commun.*, vol. 19, no. 6, pp. 4209–4219, Jun. 2020.
- [17] H. Zhou, Z. Zhang, Y. Wu, M. Dong, and V. C. M. Leung, "Energy efficient joint computation offloading and service caching for mobile edge computing: A deep reinforcement learning approach," *IEEE Trans. Green Commun. Netw.*, early access, Jul. 4, 2022, doi: 10.1109/TGCN.2022.3186403.
- [18] C. Zhong, M. C. Gursoy, and S. Velipasalar, "Deep multi-agent reinforcement learning based cooperative edge caching in wireless networks," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2019, pp. 1–6.
- [19] Y. Dai, D. Xu, K. Zhang, S. Maharjan, and Y. Zhang, "Deep reinforcement learning and permissioned blockchain for content caching in vehicular edge computing and networks," *IEEE Trans. Veh. Technol.*, vol. 69, no. 4, pp. 4312–4324, Apr. 2020.
- [20] N. Luong et al., "Applications of deep reinforcement learning in communications and networking: A survey," *IEEE Commun. Surveys Tuts.*, vol. 21, no. 4, pp. 3133–3174, 4th Quart., 2019.
- [21] X. Ye, M. Li, P. Si, R. Yang, Z. Wang, and Y. Zhang, "Collaborative and intelligent resource optimization for computing and caching in IoV with blockchain and MEC using A3C approach," *IEEE Trans. Veh. Technol.*, vol. 72, no. 2, pp. 1449–1463, Feb. 2023, doi: 10.1109/TVT.2022.3210570.
- [22] B. Jedari, G. Premsankar, G. Illahi, M. D. Francesco, A. Mehrabi, and A. Ylä-Jääski, "Video caching, analytics, and delivery at the wireless edge: A survey and future directions," *IEEE Commun. Surveys Tuts.*, vol. 23, no. 1, pp. 431–471, 1st Quart., 2021.
- [23] H. Wu and H. Lu, "Delay and power tradeoff with consideration of caching capabilities in dense wireless networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 10, pp. 5011–5025, Oct. 2019.
- [24] M. Yan, C. A. Chan, W. Li, L. Lei, A. F. Gygax, and I. Chih-Lin, "Assessing the energy consumption of proactive mobile edge caching in wireless networks," *IEEE Access*, vol. 7, pp. 104394–104404, 2019.
- [25] X. Zhang, Z. Qi, G. Min, W. Miao, Q. Fan, and Z. Ma, "Cooperative edge caching based on temporal convolutional networks," *IEEE Trans. Parallel Distrib. Syst.*, vol. 33, no. 9, pp. 2093–2105, Sep. 2022.
- [26] K. Zhang, J. Cao, S. Maharjan, and Y. Zhang, "Digital twin empowered content caching in social-aware vehicular edge networks," *IEEE Trans. Computat. Social Syst.*, vol. 9, no. 1, pp. 239–251, Feb. 2022.
- [27] J. Zhao, X. Sun, Q. Li, and X. Ma, "Edge caching and computation management for real-time Internet of Vehicles: An online and distributed approach," *IEEE Trans. Intell. Transp. Syst.*, vol. 22, no. 4, pp. 2183–2197, Apr. 2021.
- [28] W. C. Ao and K. Psounis, "Fast content delivery via distributed caching and small cell cooperation," *IEEE Trans. Mobile Comput.*, vol. 17, no. 5, pp. 1048–1061, May 2018.
- [29] Z. Gu, H. Lu, P. Hong, and Y. Zhang, "Reliability enhancement for VR delivery in mobile-edge empowered dual-connectivity sub-6 GHz and mmWave HetNets," *IEEE Trans. Wireless Commun.*, vol. 21, no. 4, pp. 2210–2226, Apr. 2022.
- [30] S. Zhang and J. Liu, "Optimal probabilistic caching in heterogeneous IoT networks," *IEEE Internet Things J.*, vol. 7, no. 4, pp. 3404–3414, Apr. 2020.
- [31] R. Wang, Z. Kan, Y. Cui, D. Wu, and Y. Zhen, "Cooperative caching strategy with content request prediction in Internet of Vehicles," *IEEE Internet Things J.*, vol. 8, no. 11, pp. 8964–8975, Jun. 2021.
- [32] 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, Jun. 2019.
- [33] 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, Nov. 2018.
- [34] G. Qiao, S. Leng, S. Maharjan, Y. Zhang, and N. Ansari, "Deep reinforcement learning for cooperative content caching in vehicular edge computing and networks," *IEEE Internet Things J.*, vol. 7, no. 1, pp. 247–257, Jan. 2020.

- [35] X. Wang, C. Wang, X. Li, V. C. M. Leung, and T. Taleb, "Federated deep reinforcement learning for Internet of Things with decentralized cooperative edge caching," *IEEE Internet Things J.*, vol. 7, no. 10, pp. 9441–9455, Oct. 2020.
- [36] H. Tian et al., "CoPace: Edge computation offloading and caching for self-driving with deep reinforcement learning," *IEEE Trans. Veh. Technol.*, vol. 70, no. 12, pp. 13281–13293, Dec. 2021.
- [37] X. Li, X. Wang, P.-J. Wan, Z. Han, and V. C. M. Leung, "Hierarchical edge caching in device-to-device aided mobile networks: Modeling, optimization, and design," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 8, pp. 1768–1785, Aug. 2018.
- [38] X. Chen, L. Jiao, W. Li, and X. Fu, "Efficient multi-user computation offloading for mobile-edge cloud computing," *IEEE/ACM Trans. Netw.*, vol. 24, no. 5, pp. 2795–2808, Oct. 2016.
- [39] J. Hu and M. P. Wellman, "Nash Q-learning for general-sum stochastic games," *J. Mach. Learn. Res.*, vol. 4, pp. 1039–1069, Nov. 2003.
- [40] P. Casgrain, B. Ning, and S. Jaimungal, "Deep Q-learning for Nash equilibria: Nash-DQN," 2019, *arXiv:1904.10554*.
- [41] Z. Zhang, D. Wang, and J. Gao, "Learning automata-based multiagent reinforcement learning for optimization of cooperative tasks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 10, pp. 4639–4652, Oct. 2021.



Shibo He (Senior Member, IEEE) is currently a Full Professor with Zhejiang University and also with the Alibaba-Zhejiang University Joint Research Institute of Frontier Technologies. His research interests include the Internet of Things, crowdsensing, and big data.

He served as the TPC Co-Chair for IEEE ScalCom 2014, the Finance and Registration Chair for ACM MobiHoc 2015, and the Symposium Co-Chair for IEEE ICC 2017 and IEEE GLOBECOM 2020. He serves on the editorial board for several journals, including IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY.

Huan Zhou (Member, IEEE) received the Ph.D. degree from the Department of Control Science and Engineering, Zhejiang University. He was a Visiting Scholar with Temple University from November 2012 to May 2013. He was a CSC supported Post-Doctoral Fellow with The University of British Columbia from November 2016 to November 2017. He is currently a Professor with the College of Computer, Northwestern Polytechnical University. He has published more than 70 research papers in some international journals and conferences, including IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS, IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, IEEE TRANSACTIONS ON MOBILE COMPUTING, and IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS. His research interests include mobile social networks, VANETs, opportunistic mobile networks, and mobile data offloading. He was a TPC Member of IEEE MASS, ICCCN, GLOBECOM, ICC, and WCSP. He received the Best Paper Award of I-SPAN 2014 and I-SPAN 2018. He was the TPC Chair of EAI BDTA 2020 and the Local Arrangement Chair of I-SPAN 2018. He was a Lead Guest Editor of *Pervasive and Mobile Computing*. He is serving as an Associate Editor for EURASIP Journal on Wireless Communications and Networking and IEEE ACCESS.



Geyong Min (Senior Member, IEEE) received the B.Sc. degree in computer science from the Huazhong University of Science and Technology, China, in 1995, and the Ph.D. degree in computing science from the University of Glasgow, U.K., in 2003. He is currently a Professor of high performance computing and networking with the Department of Computer Science, College of Engineering, Mathematics, and Physical Sciences, University of Exeter, U.K. His research interests include computer networks, wireless communications, parallel and distributed computing, ubiquitous computing, multimedia systems, and modeling and performance engineering.

Kai Jiang (Student Member, IEEE) is currently pursuing the Ph.D. degree in cyberspace security with Wuhan University, China. His research interests include mobile edge computing, deep reinforcement learning, and the Internet of Vehicles.



Jie Wu (Fellow, IEEE) is the Chair and a Laura H. Carnell Professor with the Department of Computer and Information Sciences, Temple University, Philadelphia, PA, USA. Prior to joining Temple University, he was the Program Director of the National Science Foundation and a Distinguished Professor with Florida Atlantic University. He has regularly published in scholarly journals, conference proceedings, and books. His current research interests include mobile computing and wireless networks, routing protocols, cloud and green computing, network trust and security, and social network applications. He was a recipient of the 2011 China Computer Federation (CCF) Overseas Outstanding Achievement Award. He was the General Co-Chair/Chair of IEEE MASS 2006, IEEE IPDPS 2008, IEEE ICDCS 2013, ACM MobiHoc 2014, IEEE ICPP 2016, and IEEE CNS 2016, as well as the Program Co-Chair of IEEE INFOCOM 2011 and CCF CNCC 2013. He was the Chair of the IEEE Technical Committee on Distributed Processing (TCDP). He serves on several editorial boards, including IEEE TRANSACTIONS ON SERVICES COMPUTING and the *Journal of Parallel and Distributed Computing*. He was an IEEE Computer Society Distinguished Visitor and an ACM Distinguished Speaker. He is a CCF Distinguished Speaker.

1228	[35]	X. Wang, C. Wang, X. Li, V. C. M. Leung, and T. Taleb, "Federated deep reinforcement learning for Internet of Things with decentralized cooperative edge caching," <i>IEEE Internet Things J.</i> , vol. 7, no. 10, pp. 9441–9455, Oct. 2020.	1278
1229	[36]	H. Tian et al., "CoPace: Edge computation offloading and caching for self-driving with deep reinforcement learning," <i>IEEE Trans. Veh. Technol.</i> , vol. 70, no. 12, pp. 13281–13293, Dec. 2021.	1279
1230	[37]	X. Li, X. Wang, P.-J. Wan, Z. Han, and V. C. M. Leung, "Hierarchical edge caching in device-to-device aided mobile networks: Modeling, optimization, and design," <i>IEEE J. Sel. Areas Commun.</i> , vol. 36, no. 8, pp. 1768–1785, Aug. 2018.	1280
1231	[38]	X. Chen, L. Jiao, W. Li, and X. Fu, "Efficient multi-user computation offloading for mobile-edge cloud computing," <i>IEEE/ACM Trans. Netw.</i> , vol. 24, no. 5, pp. 2795–2808, Oct. 2016.	1281
1232	[39]	J. Hu and M. P. Wellman, "Nash Q-learning for general-sum stochastic games," <i>J. Mach. Learn. Res.</i> , vol. 4, pp. 1039–1069, Nov. 2003.	1282
1233	[40]	P. Casgrain, B. Ning, and S. Jaimungal, "Deep Q-learning for Nash equilibria: Nash-DQN," 2019, <i>arXiv:1904.10554</i> .	1283
1234	[41]	Z. Zhang, D. Wang, and J. Gao, "Learning automata-based multiagent reinforcement learning for optimization of cooperative tasks," <i>IEEE Trans. Neural Netw. Learn. Syst.</i> , vol. 32, no. 10, pp. 4639–4652, Oct. 2021.	1284
1235			1285
1236			1286
1237			1287
1238			1288
1239			1289
1240			1290
1241			1291
1242			1292
1243			1293
1244			1294
1245			1295
1246			1296
1247			1297
1248			1298
1249			1299
1250			1300
1251			1301
1252			1302
1253			1303
1254			1304
1255			1305
1256			1306
1257			1307
1258			1308
1259			1309
1260			1310
1261			1311
1262			1312
1263			1313
1264			1314
1265			1315
1266			1316
1267			1317
1268			1318
1269			1319
1270			1320
1271			1321
1272			1322
1273			1323
1274			1324

AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES

PLEASE NOTE: We cannot accept new source files as corrections for your article. If possible, please annotate the PDF proof we have sent you with your corrections and upload it via the Author Gateway. Alternatively, you may send us your corrections in list format. You may also upload revised graphics via the Author Gateway.

Carefully check the page proofs (and coordinate with all authors); additional changes or updates WILL NOT be accepted after the article is published online/print in its final form. Please check author names and affiliations, funding, as well as the overall article for any errors prior to sending in your author proof corrections. Your article has been peer reviewed, accepted as final, and sent in to IEEE. No text changes have been made to the main part of the article as dictated by the editorial level of service for your publication.

AQ:1 = Please confirm or add details for any funding or financial support for the research of this article.

AQ:2 = Note that if you require corrections/changes to tables or figures, you must supply the revised files, as these items are not edited for you.

Distributed Deep Multi-Agent Reinforcement Learning for Cooperative Edge Caching in Internet-of-Vehicles

Huan Zhou[✉], Member, IEEE, Kai Jiang[✉], Student Member, IEEE, Shibo He[✉], Senior Member, IEEE,
Geyong Min[✉], Senior Member, IEEE, and Jie Wu[✉], Fellow, IEEE

Abstract—Edge caching is a promising approach to reduce duplicate content transmission in Internet-of-Vehicles (IoVs). Several Reinforcement Learning (RL) based edge caching methods have been proposed to improve the resource utilization and reduce the backhaul traffic load. However, they only obtain the local sub-optimal solution, as they neglect the influence from environments by other agents. This paper investigates the edge caching strategies with consideration of the content delivery and cache replacement by exploiting the distributed Multi-Agent Reinforcement Learning (MARL). A hierarchical edge caching architecture for IoVs is proposed and the corresponding problem is formulated with the goal to minimize the long-term content access cost in the system. Then, we extend the Markov Decision Process (MDP) in the single agent RL to the context of a multi-agent system, and tackle the corresponding combinatorial multi-armed bandit problem based on the framework of a stochastic game. Specifically, we firstly propose a Distributed MARL-based Edge caching method (DMRE), where each agent can adaptively learn its best behaviour in conjunction with other agents for intelligent caching. Meanwhile, we attempt to reduce the computation complexity of DMRE by parameter approximation, which legitimately simplifies the training targets. However, DMRE is enabled to represent and update the parameter by creating a lookup table, essentially a tabular-based method, which generally performs inefficiently in large-scale scenarios. To circumvent the issue and make more expressive parametric models, we incorporate the technical advantage of the Deep-Q Network into DMRE, and further develop a computationally efficient method (DeepDMRE) with neural network-based Nash

equilibria approximation. Extensive simulations are conducted to verify the effectiveness of the proposed methods. Especially, DeepDMRE outperforms DMRE, *Q*-learning, LFU, and LRU, and the edge hit rate is improved by roughly 5%, 19%, 40%, and 35%, respectively, when the cache capacity reaches 1,000 MB.

Index Terms—Edge caching, Internet-of-Vehicles, content delivery, cache replacement, multi-agent reinforcement learning.

I. INTRODUCTION

WITH the explosive growth of vehicles on the road and the progress of wireless multi-access technology, we have witnessed unprecedented changes in the traditional transportation system, which has evolved from a technology-driven era to a more powerful data-driven intelligent era [1]. As a fundamental paradigm of 5G networks, Internet-of-Vehicles (IoVs) enables a wide variety of vehicles to provide reliable vehicular multimedia services, cooperative cruise control, and path navigation, etc., which have paved a path towards intelligent transportation, and improved the driving safety and travel comfort of passengers [2], [3].

Meanwhile, these flourishing vehicular applications require vehicles to access vast amounts of Internet data, especially for some delay-sensitive contents, leading to severe network congestion and considerable delay in content delivery [4], [5], [6]. Technically, tight Quality-of-Service (QoS) requirements are generally placed for such applications, which makes cloud-based processing architectures infeasible due to the long transmission distance and limited backhaul link capacity inherent in content delivery from remote cloud servers. Indeed, the rapid increase in delay and unreliability driven by the transmission distance has become an inevitable issue for supporting massive content delivery and gained widespread attention lately [7], [8]. Despite the continuous efforts made on improving the backhaul link capacity through advanced technologies, the radio spectrum utilization has obviously reached the theoretical bound [9], [10], [11]. Therefore, these efforts are insufficient to cope with these great challenges. A fundamental innovation that breaks through the bottleneck of massive content delivery in IoVs there is urgently required.

Fortunately, large-scale data analysis shows that different contents often require different priorities. Only a few popular contents account for the majority of downloads, while the access demand for most of the rest is relatively

Manuscript received 15 January 2022; revised 18 July 2022, 28 October 2022, and 1 March 2023; accepted 17 April 2023. This work was supported in part by the National Natural Science Foundation of China under Grant 62172255 and Grant U1909207 and in part by the Outstanding Youth Program of Hubei Natural Science Foundation under Grant 2022CFA080. The associate editor coordinating the review of this article and approving it for publication was A. Schmeink. (*Corresponding author: Kai Jiang*.)

Huan Zhou is with the School of Computer Science, Northwestern Polytechnical University, Xi'an 710129, China, and also with the College of Computer and Information Technology, China Three Gorges University, Yichang 443002, China (e-mail: zhouchuan17@gmail.com).

Kai Jiang is with the School of Cyber Science and Engineering, Wuhan University, Wuhan 430000, China (e-mail: kai.jiang@whu.edu.cn).

Shibo He is with the College of Control Science and Technology, Zhejiang University, Hangzhou 310027, China (e-mail: s18he@zju.edu.cn).

Geyong Min is with the Department of Computer Science, College of Engineering, Mathematics, and Physical Sciences, University of Exeter, EX4 4QF Exeter, U.K. (e-mail: g.min@exeter.ac.uk).

Jie Wu is with the Department of Computer and Information Sciences, Temple University, Philadelphia, PA 19122 USA (e-mail: jiewu@temple.edu).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TWC.2023.3272348>.

Digital Object Identifier 10.1109/TWC.2023.3272348

72 small [12], [13]. This request pattern promotes the imple-
 73 mentation of edge caching technology in IoVs. Expressly,
 74 edge caching has provided an alternative to alleviate the
 75 backhaul link strain and content access delay inside the
 76 future IoVs, which pre-caches frequently-used contents close
 77 to vehicles by pushing cloud functions to intermediate Road-
 78 side Units (RSUs). Hence, edge caching enables vehicles to
 79 access popular content from the caching-enabled nearby RSUs
 80 directly, instead of repeatedly downloading from remote cloud
 81 servers [14]. In this way, the redundant traffic and transmission
 82 resource consumption at both backhaul and core networks can
 83 be significantly reduced, and the QoS is improved as well.
 84 In addition, enabling edge caching in IoVs can benefit extra
 85 technical superiorities from its distributed architectures and
 86 small-scale nature, including privacy protection, scalability,
 87 context awareness, and so on.

88 Compared with the cloud server, the cache capacity of a
 89 single edge node is insufficient to support all popular con-
 90 tents. Without edge cooperation, the overall cache utilization
 91 and effectiveness are prone to be under-utilized. Specifically,
 92 adjacent RSUs can share data traffic and exchange popular
 93 contents rather than operating individually, so the bilateral
 94 synergy between their caches and the centralized cloud should
 95 also be leveraged to facilitate good performance, which largely
 96 depends on a well-designed cooperative edge caching strategy.
 97 On the other hand, with the revival of artificial intelligence,
 98 Reinforcement Learning (RL) has exhibited its particular
 99 potential for the method design of efficient edge caching in
 100 IoVs. In RL, the agent can well capture the hidden dynamics
 101 of the environments and imitate the best features by trial-
 102 and-error interaction, thus learning a series of intelligence
 103 behaviors to enhance the edge cache utilization [15], [16], [17].
 104 However, massively diverse, highly vibrant and distributed
 105 edge caching contexts make most previous RL work unable
 106 to adapt as they neglected the environment's influence by
 107 other agents when an agent interacts and learns the settings
 108 independently. Therefore, only the local sub-optimal solution
 109 can be obtained in the system, not the optimal global one,
 110 especially in a long-term optimization problem [18].

111 Indeed, it is more reasonable and helpful to investigate from
 112 a distributed perspective under the multi-agent context, which
 113 leads to inherent robustness and high scalability.¹ Meanwhile,
 114 several studies [19], [20], and [21] have proved that agents
 115 considering the impact of joint actions perform better than
 116 corresponding agents learning only based on their own actions.
 117 Therefore, this paper investigates the edge caching strategy
 118 with consideration of the content delivery and cache replace-
 119 ment by exploiting the distributed Multi-Agent Reinforcement
 120 Learning (MARL). A hierarchical edge caching architecture
 121 for IoVs is first proposed, where the cooperative caching
 122 among multi-RSUs and Macro Base Station (MBS) is utilized
 123 to reduce the content access cost and the backhaul traffic in
 124 the system. Then, the Markov Decision Process (MDP) in
 125 the single-agent RL is extended to the context of multi-agent

¹Robustness: when one or more agents fail in the system, the remaining agents can take over some of their tasks.

Scalability: by sure design, most multi-agent systems allow easy insertion of new agents.

system and the corresponding combinatorial multi-armed bandit problem is tackled based on the framework of stochastic games. Specifically, a **Distributed MARL-based Edge caching method (DMRE)** is proposed firstly, where each agent can adaptively learn its best behaviour in conjunction with other agents for intelligent caching. Meanwhile, we attempt to reduce the computation complexity of DMRE by parameter approximation, which legitimately simplifies the training targets. However, DMRE is enabled to represent and update the parameter by creating a lookup table, essentially a tabular-based method, which generally performs inefficiently in large-scale scenarios. To circumvent the issue and make more expressive parametric models, we incorporate the technical advantage of the Deep-*Q* Network into DMRE, and further develop a computationally efficient method (**DeepDMRE**) with neural network-based Nash equilibria approximation. The main contributions are summarized as follows:

- 1) An original hierarchical edge caching architecture for IoVs is presented and the corresponding problem is formulated with the goal to minimize the long-term content access cost in the system.
- 2) To address the overhead minimization problem, an MDP for the optimization of cache replacement process in an available RSU is defined. Then, the MDP is extended to the context of a multi-agent system, and the corresponding combinatorial multi-armed bandit problem is tackled with the proposed DMRE firstly.
- 3) Furthermore, to circumvent the issue by DMRE and make more expressive parametric models, we incorporate the technical advantage of the Deep-*Q* Network into DMRE, and further develop a computationally efficient method (DeepDMRE) with neural network-based Nash equilibria approximation.
- 4) Extensive simulation results demonstrate that DeepDMRE significantly outperforms other benchmark methods in different scenarios. Meanwhile, relevant theoretical proofs of convergence and feasibility of the proposed methods are presented at the end of the paper.

The remainder of this paper is organized as follows. Section II and Section III describe the related work and the system framework, respectively. Section IV formulates the problem to minimize the long-term content access cost in the system. Section V extends MDP to a multi-agent system under the stochastic game framework, and proposes two MARL-based methods. Moreover, simulation results are analyzed in Section VI. Conclusion is summarized in Section VII.

II. RELATED WORK

Recently, extensive works have focused on content caching strategies at the edge of networks [22], [23], [24], [25], [26], [27], [28], [29], [30]. Jedari et al. [22] elaborated an exhaustive review on edge caching and discussed its implementation and outlook. Wu and Lu [23] investigated the tradeoff between content delivery delay and power consumption in small-cell network caching, and proposed an iterative algorithm to approach the optimal tradeoff between these two metrics. Yan et al. [24] exploited the energy assessment

models for mobile edge caching and proposed a predictive caching strategy to potentially improve the cache hit rates. Zhang et al. [25] focused on a learning-based edge caching scheme to enable cooperation among different edge nodes with limited resources, intending to reduce the content delivery latency and maximize the overall content caching value. Similarly, to maximize the utility of the caching system, a vehicular edge caching mechanism based on user preference similarity and service availability was proposed in [26]. Furthermore, Zhao et al. [27] dynamically orchestrated the cache resources in IoVs by leveraging the combined power of Lyapunov optimization, matching theory, and consensus alternating direction method of multipliers. In [28], Ao et al. leveraged the idea of coordinated multi-point and devised an efficient cooperative caching strategy to maximize the system throughput, where the available cache size and the network topology are both taken into account. In [29], Gu et al. proposed a dual-connectivity sub-6 GHz and mmWave heterogeneous network architecture to conduct a collaborative edge resource design, aiming to maximise virtual reality delivery reliability. In [30], Zhang et al. considered the stochastic geometry based probabilistic caching to improve the offloading rate for backhaul links, and presented an improved caching probability conversion algorithm to obtain the closed-form solutions. These works prove the advantage of edge caching. However, most of them are quasi-static and myopic, and thus perform unsatisfactorily in the diverse and dynamic IoVs system.

As Reinforcement Learning (RL) can well capture the hidden dynamics of the environments, it has been exploited to enhance the intelligence of IoVs. Specifically, Wang et al. [31] clustered vehicles to simplify the process of content requesting via the K-means method, and proposed an RL-based cooperative caching strategy with content request prediction in IoVs. In [32], Dai et al. developed a novel edge caching architecture for AI-empowered vehicular networks. They formulated a joint optimization of edge computing and caching to maximize the system utility, and exploited an RL based method to dynamically allocate computing and caching resources. Without any information or assumptions about the environment, Tan and Hu [33] achieved an RL based adaptive framework to tackle the resource allocation of communication, caching, and computing in vehicular networks. This framework was also used by Qiao et al. [34] to study the problem of content placement and content delivery in IoVs, which aims to reduce the system cost under the constraint of content delivery latency. Furthermore, Wang et al. [35] proposed a federated deep RL-based cooperative edge caching framework to eliminate duplicate traffic and improve the edge cache utilization with specific QoS requirements. Tian et al. [36] considered the mobility and changing requests of self-driving vehicles, and proposed a deep RL-based collaborative caching method to break the curse of high dimensionality in the state space of MDP, as well as decrease the redundant transmission delay.

The aforementioned studies are based on centralized methods that may not always be feasible in practice due to distributed and cooperative network topology. The fundamental problem is that these methods treat the other agents as a part of the environment and thus neglect the influence on the

TABLE I
NOTATIONS AND SYMBOLS

Notation	Explanation
\mathcal{N}	The set of all RSUs
\mathcal{F}	The index set of available contents
s_f	The size of content f
G_i	The limited available cache capacity of RSU i
\mathcal{U}	The set of all requested vehicles
ρ_f	The content popularity of content f
$T_{u,i}$	The content transmission delay between vehicle u and RSU i
$T_{i,j}$	The content transmission delay between RSU i and RSU j
$T_{i,N+1}$	The content transmission delay between RSU i and the MBS
$r_{u,i}$	The wireless downlink transmission rate between vehicle u and RSU i
$C_{u,i}$	The incurred transmission cost between vehicle u and RSU i
$C_{i,j}$	The incurred transmission cost between RSU i and RSU j
$C_{i,N+1}$	The incurred transmission cost between RSU i and the MBS
$\mathcal{L}_{f,i,j}$	The total overhead through different links
$a_{f,i,j}^t$	The caching state of local RSU i for the requested content f
$c_{f,i}^t$	The content replacement decision of RSU i in time slot t
$q_{i,\mathcal{F}}^t$	The content request arrived from the vehicles in time slot t

AQ:2

environment by other agents' behavior. Different from them, this paper investigates the edge caching strategy considering the content delivery and cache replacement in a distributed perspective. Specifically, the DMRE framwork is proposed to reduce the backhaul traffic in the system, where each agent can adaptively learn its best behavior in conjunction with other agents for intelligent caching. In addition, to further circumvent the issue by DMRE and make more expressive parametric models, another computationally efficient method (DeepDMRE) is also proposed, which approximates the Nash equilibria by constructing neural networks.

III. DISTRIBUTED VEHICULAR EDGE NETWORK ARCHITECTURE

This section presents the cooperative edge caching-supported IoVs architecture, including network infrastructure, content popularity pattern, as well as the process of content delivery and cache replacement. The main notations with their descriptions are listed in Table I.

A. Network Architecture

In this paper, a cooperative edge caching-supported IoVs architecture with an MBS and N small cells is considered, each of which is naturally equipped with an RSU, as shown in Fig. 1. The set of these RSUs is denoted by $\mathcal{N} = \{1, 2, \dots, N\}$. Considering the abundant cache capacity, it is assumed that MBS (denoted by $N + 1$ in the system) is connected to the Service Providers (SPs) and can cache all contents. Here, let $\mathcal{F} = \{1, 2, \dots, F\}$ represent the content library with total F available contents that are supported by the SPs, and s_f is the size of content f ($f \in \mathcal{F}$). Besides, each RSU i ($i \in \mathcal{N}$) is endowed with a limited available cache capacity of G_i , such that a small portion of popular contents can be cached in the RSUs to eliminate the redundant traffic.

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

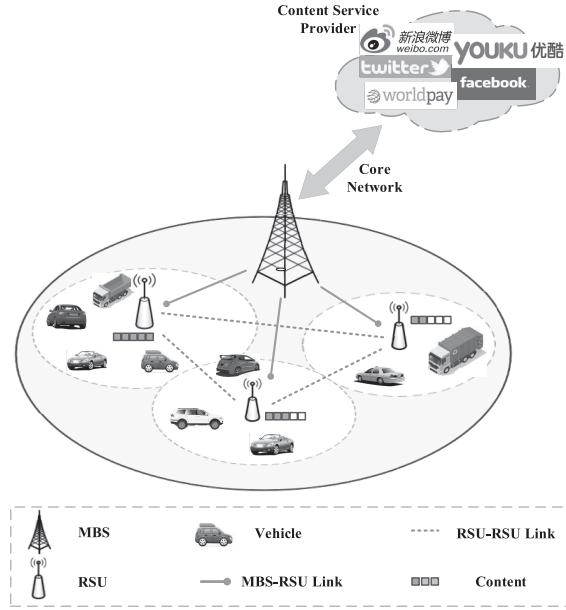


Fig. 1. Cooperative Edge Caching-supported IoVs Architecture.

All RSUs in a cluster form are connected to the remote MBS via wired lines. U vehicles (denoted as $\mathcal{U} = \{1, 2, \dots, U\}$) are randomly scattered in the wireless coverage scope and request various contents frequently via cellular links. Overlapping coverage between neighbouring cells is not considered, and these vehicles are assumed to be located in at least one small cell and will be served by the RSU therein; the set of vehicles in a small cell may change dynamically due to the vehicle mobility. Each RSU serves the local content requests within its coverage range. Meanwhile, the system is assumed to operate in a fixed length of time slots (i.e., time sequence) $t \in \{1, \dots, \Gamma\}$, where Γ denotes the finite time horizon. In each time slot, a vehicle can only request one content, while the caching decision of each RSU is also updated periodically.

As vehicles may exhibit different preferences for popular content, we introduce content popularity $\rho_f(f \in \mathcal{F})$ to reflect the content access requirements in the content library, i.e., ρ_f denotes the conditional probability of requesting content f among various vehicles' content requests. Inspired by existing works [27], [37], the content popularity is assumed to obey the Mandelbrot-Zipf distribution here. Then, the expected requesting probability for content f is obtained as:

$$\rho_f = \frac{(R_f + \tau)^{-\vartheta}}{\sum_{i \in \mathcal{F}} (R_i + \tau)^{-\vartheta}}, \quad \forall f \in \mathcal{F}, \quad (1)$$

where R_f is the rank of content f in the descending order of content popularity; $\vartheta > 0$ and τ are the skewness factor and plateau factor characterizing the distribution, respectively. Notably, due to the vehicle mobility, the local content popularity can be highly spatio-temporally varying over time slots and the requested vehicles in a cell, making caching strategy design more challenging.

Within the coverage of MBS, each RSU is able to cache various contents to meet the content access requirements from vehicles. Caching contents in RSUs enables the content

request to be accommodated in proximity, without duplicate downloading from remote MBS. Particularly, each RSU can determine where the content request is answered and which local cache should be replaced. Once a vehicle sends an access request within the range of the small cell, the local RSU will traverse its cache space firstly to check whether the desired content is cached in itself. If the content is not sought in its cache, the local RSU can either retrieve the content from the adjacent RSUs or download the content directly from the MBS, and deliver it to the requested vehicle later. Meanwhile, all RSUs may replace their cache with popular contents according to the content requests of each time slot. Underlying these is a sequential decision-making problem, which corresponds to making efficient content delivery and cache replacement decisions under constraints.

B. Content Delivery Model

$T_{u,i}$, $T_{i,j}$, and $T_{i,N+1}$ are used to denote the content transmission delay between vehicle u and RSU i , RSU i and RSU j , RSU i and MBS, respectively. It is worth noticing that the transmission delay of sending request identifier by the vehicle is neglected due to its tiny data size and the high link rate in most instances. Specifically, transmission delay of content f from any requested vehicle u ($u \in \mathcal{U}$) to its connected local RSU i can be calculated as $T_{u,i} = s_f / r_{u,i}$, where $r_{u,i}$ is the wireless downlink data rate between vehicle u and RSU i . Here, we assume that each RSU consists of multiple channels, and the allocated bandwidth to each channel is the same [38]. Therefore, considering the large-scale fading, the wireless downlink data rate $r_{u,i}$ is derived by

$$r_{u,i} = B_{u,i} \log_2 \left(1 + \frac{P_i g_{u,i}}{\sigma^2 + \sum_{v \in \mathcal{U} \setminus \{u\}} P_v g_{v,i}} \right), \quad (2)$$

where $B_{u,i}$ denotes the channel bandwidth, P_i denotes the transmission power of RSU i , σ^2 and $g_{u,i}$ are the white Gaussian noise power and the wireless propagation channel gain, respectively. Interference management is also considered for this wireless communication scenario. Furthermore, since the communications of MBS-RSU and RSU-RSU are via wired optical cables, the transmission delay $T_{i,N+1}$ and $T_{i,j}$ are set to constant values [35].

Delivering desired contents to vehicles will introduce extra transmission costs for caching nodes (cooperative RSUs or the MBS, denoted as $\mathcal{H} = \mathcal{N} \cup \{N+1\}$). As illustrated in Fig. 2, three components of transmission cost exist due to different cache hit situations. Here, $C_{u,i}$, $C_{i,j}$, and $C_{i,N+1}$ are used to denote the incurred transmission cost through RSU-Vehicle, RSU-RSU, and MBS-RSU links, respectively. Specifically, the incurred transmission cost through the RSU-Vehicle link is equal to the product of the link bandwidth, the content delivery delay, and the price of the RSU-Vehicle link's unit leased communication resource, which is expressed as:

$$C_{u,i} = B_{u,i} T_{u,i} \xi_{u,i}, \quad (3)$$

where $\xi_{u,i}$ represents the price of unit leased communication resource of wireless link. Similarly, the incurred transmission cost through RSU-RSU and MBS-RSU links can also be obtained by this way.

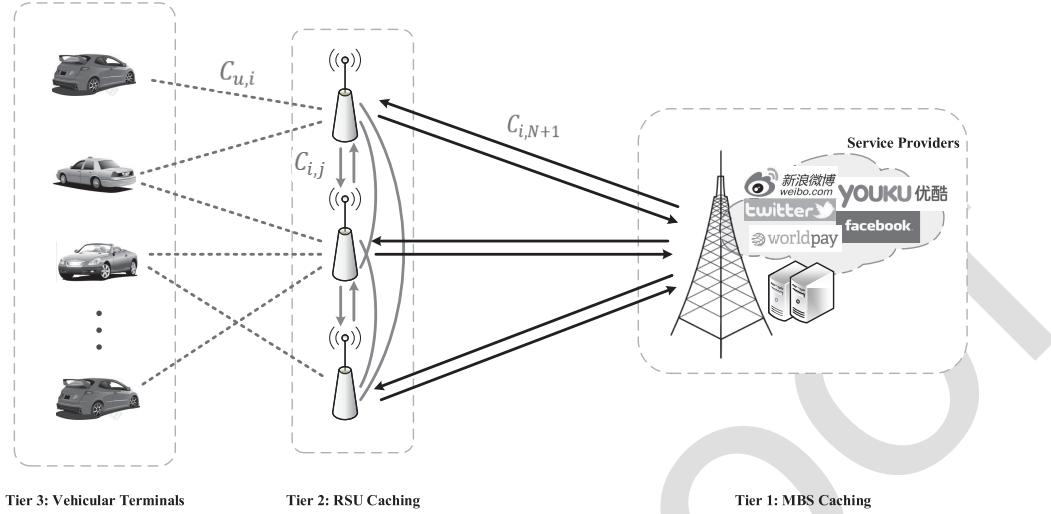


Fig. 2. Content delivery process and aggregate overhead in hierarchical vehicle edge networks.

Here, the binary decision variable $a_{f,i,j}^t \in \{0, 1\}$, $j \in \mathcal{H}$ is used to denote the caching state of any local RSU i for the requested content f in a given time slot t . When a content request occurs, the aggregate overhead of different links can be analyzed as follows:

- $a_{f,i,i}^t = 1$ indicates that content f is cached in the local RSU i at time slot t and can be delivered to the vehicle directly. In this case, the total overhead $\mathcal{L}_{f,i,i}$ is just the transmission cost $C_{u,i}$ from the local RSU to the vehicle;
- $a_{f,i,j}^t = 1$ ($j \in \mathcal{N} \setminus \{i\}$) indicates that content f is not cached in the local RSU i , but at least one of RSUs in the system has cached the desired content, thus the local RSU can inquire and obtain it from the cooperative RSU j . Then, the total overhead $\mathcal{L}_{f,i,j}$ consists of the transmission cost $C_{i,j}$ from the cooperative RSU j to the local RSU i and the transmission cost $C_{u,i}$;
- $a_{f,i,N+1}^t = 1$ indicates that there is no expected content in all cooperative RSUs, and the local RSU i has to forward the request identifier to MBS for processing at time slot t , i.e., the local RSU i downloads content f from MBS directly. In this way, the total overhead $\mathcal{L}_{f,i,N+1}$ consists of the transmission cost $C_{i,N+1}$ from MBS to the local RSU i and the transmission cost $C_{u,i}$.

C. Cache Replacement Model

In a diverse and dynamic edge caching scenario, it is indispensable to elaborate a wise content replacement strategy for efficiently managing the cache space. In order to meet the content requests that arrive later, each cooperative RSU should update the less popular contents in each operation phase to further improve the utilization and effectiveness of the whole cache.

Generally, the local RSU needs to determine whether to keep or replace its cache contents. When the newly requested contents arrive from other adjacent RSUs or MBS, some contents in the local RSU's existing cache may be evicted because of the finite cache space. As a result, the problem is whether and which contents should be replaced with the

new ones when the cache space is fully occupied. To this end, we represent the replacement control of RSU i by $c_{\mathcal{F},i}^t = \{c_{1,i}^t, \dots, c_{f,i}^t, \dots, c_{F,i}^t\}$ in time slot t , where $c_{f,i}^t \in \{0, 1\}$ ($f \in \mathcal{F}$) indicates whether and which content in RSU i should be replaced by the current content, e.g., $c_{f,i}^t = 1$ means that content f in RSU i needs to be replaced by the current arrived one, while $c_{f,i}^t = 0$ is the opposite. To efficiently store contents in RSUs, a utility-based replacement strategy will be introduced below.

IV. PROBLEM FORMULATION

This section formulates the optimization problem for IoV systems, and models the cache replacement process in an RSU as an MDP.

A. Objective Function

The RSUs in a cluster form can explore the potential of cooperation to fully utilize the cache resources. Considering the anticipated future utility, we intend to avoid myopic caching decisions while minimizing the long-term content access cost in the system. The corresponding problem is formulated as follows:

$$\begin{aligned} & \min_{a_{f,i,j}} \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \sum_{f=1}^F \sum_{i=1}^N \sum_{j=1}^{N+1} \rho_f a_{f,i,j}^t \mathcal{L}_{f,i,j}^t \\ & \text{s. t.} \quad C1 : a_{f,i,j}^t \in \{0, 1\}, \quad \forall f, \forall i, \forall j \\ & \quad C2 : a_{f,N+1,N+1}^t = 1, \quad \forall f \\ & \quad C3 : a_{f,i,j}^t \leq a_{f,j,j}^t, \quad \forall f, \forall i, \forall j \\ & \quad C4 : \sum_{j=1}^{N+1} a_{f,i,j}^t = 1, \quad \forall f, \forall i \\ & \quad C5 : \sum_{f=1}^F a_{f,i,i}^t s_f \leq G_i, \quad \forall i \end{aligned} \quad (4)$$

Here, $\lim_{T \rightarrow \infty} \frac{1}{T} (\dots)$ is the time averaged overhead of content delivery in the system. The meaning of the above constraints is

as follows: C1 guarantees the constraint of the binary caching decision's integer nature; C2 shows the MBS's sufficient cache capacity and guarantees that MBS has cached all available contents; C3 denotes that only the cooperative RSUs or MBS, which have cached the related content, can answer the content request; C4 ensures that the content requested by a particular vehicle can only be served by an RSU or MBS ultimately; C5 is used to promise that the cache usage of each RSU should be less than its cache capacity.

433 B. Problem and Challenges

434 In fact, it is challenging to solve the above problem
 435 in Eq. (4) directly since we have to obtain the optimal
 436 coordination of the caching decision variables $a_{f,i,j}^t$ ($f \in$
 $\mathcal{F}, j \in \mathcal{H} \setminus \{i\}$) in each time slot. However, the caching deci-
 437 sion variable $a_{f,i,j}^t$ of any RSU is binary and time-varying,
 438 which aggravates the difficulty by solving the problem at a
 439 time slot exhaustively. The system has to collect multitudinous
 440 network state information and make the global decision on
 441 cache replacement and content delivery for each RSU. More-
 442 over, we are devoted to a more practical case in which the prior
 443 information of the content request pattern is unknown. Thus,
 444 the optimization of the system is a combinatorial multi-armed
 445 bandit problem, and the objective function is undoubtedly
 446 NP-hard. Since the feasible set of the problem is not convex
 447 and the complexity is enormous, conventional methods using
 448 model-based heuristics or predefined rules may be unadaptable
 449 to make intelligent caching decisions under the dynamic
 450 system property.

452 C. Markov Decision Process

453 The MDP model is typically used to describe almost all
 454 sequential decision-making problems. Here, we model the
 455 optimization of cache replacement process in an RSU as an
 456 MDP. Three critical elements are identified as follows:

- **State Space:** The state in MDP is a space to reflect the IoV's environment. Accordingly, the state of an available RSU i is determined by the realization of the current caching situations and the request demand situations, which can be given as $z_t^i = \{a_{\mathcal{F},i,i}^t, q_{i,\mathcal{F}}^t\}$. As described earlier, the former $a_{\mathcal{F},i,i}^t = \{a_{1,i,i}^t, \dots, a_{f,i,i}^t, \dots, a_{F,i,i}^t\}$ denotes the current caching state respecting to the contents in RSU i . The latter $q_{i,\mathcal{F}}^t = \{q_{i,1}^t, \dots, q_{i,f}^t, \dots, q_{i,F}^t\}$ is the arrived content request state from vehicles in time slot t . Specifically, $q_{i,f}^t = 1$ means that at least one UE within local RSU i requests for content f in time slot t , $q_{i,f}^t = 1$ is the opposite. The above information will be aggregated as an input state in each time slot.
- **Action Space:** Once receiving a state in each time slot, the agent is responsible for deciding the caching and delivery strategy. Therefore, with the current state z_t^i , the action d_t^i involves two parts, $a_{\mathcal{F},i,j}^t$ and $c_{\mathcal{F},i}^t$, where matrix $a_{\mathcal{F},i,j}^t$ encodes the decision of RSU i for the requested contents, and the vector $c_{\mathcal{F},i}^t$ indicates the content replacement control in RSU i in slot t .

- **Reward Function:** Defining an appropriate reward function is vital since the reward value is the metric for each agent to evaluate the action quality. Generally speaking, the objective function is related to the immediate reward. Our objective in the related problem is to minimize the long-term content access cost in the system, which is the inverse goal of an agent that tries to achieve the maximum cumulative discounted reward. Thus, the reward function should be negatively correlated with the optimization objective. For this purpose, the immediate reward is defined as normalized $r(z_t^i, d_t^i) = e^{-\sum_f^F \sum_j^{N+1} \rho_f a_{f,i,j}^t L_{f,i,j}^t}$, where we utilize a negative exponential function to transform the problem legitimately. Furthermore, the reward value of an agent should satisfy the constraint conditions to ensure the validity of the results. Conversely, punishment negative will be incorporated into the reward if constraints are violated in Eq. (4).

V. MARL-EMPOWERED COOPERATIVE EDGE CACHING

Compared with dynamic programming methods, it is more suitable to adopt RL-based methods when the transition probability is ambiguous in MDP. This section extends MDP to a multi-agent system under the framework of a stochastic game, and proposes two MARL-based edge caching methods.

A. Markov Game Model

According to the MDP above, we have proposed a Q -learning based edge caching method in our previous work [7]. Q -learning is a model-free independent RL method based on value iteration. Each agent in Q -learning learns the settings separately through repeated interaction, and regards other agents as a part of the environment. Although single-agent properties are transferred directly to multi-agent contexts in Q -learning based method, the IoVs system is formulated as a multitude of subsystems where the cache replacement process in each available RSU is optimized individually. Nevertheless, as an independent caching node, each RSU should have its own caching strategy according to the corresponding content request and caching state, which may be very different between different RSUs. In addition, the caching strategies of agents are mutually influential in the distributed edge caching scenarios. Q -learning based method and their variants can only obtain the local sub-optimal solution, as they do not consider the influence on the environment by other agents when an agent interacts separately, i.e., each agent greedily makes decisions on its own benefit.

As far as these issues are concerned, we consider extending the MDP to a multi-agent system and formulating the cache replacement process as a Markov (a.k.a. Stochastic) Game (MG) model innovatively. In the MG model, due to the interaction between agents, the optimal decision of a single agent cannot ensure the optimal global solution of the system. Agents should not only observe their own immediate rewards and actions taken previously, but also those of others as well. Our objective is to obtain the best strategy for every agent. Therefore, it is necessary to make joint decisions among all

agents to avoid the deviation in a single decision, so as to learn the coordination and optimization of the whole system.

In MG model, agents execute actions at the same time, and the rewards of agents are generally arbitrarily related. Tuple $\{N, \mathcal{Z}, D_1, \dots, D_N, p, r_1, \dots, r_N, \gamma\}$ can characterize a standard formal definition of MG model, where N is the number of agents; \mathcal{Z} is the finite state space of the environment, each state $z_{t+1} \in \mathcal{Z}$ in the game consists of the individual states of all agents in the system; D_i stands for a discrete action space of agent i ($i = 1, \dots, N$); the joint action space of all agents is represented as $\mathcal{D} = D_1 \times \dots \times D_N$, and for notational convenience, we use $\vec{d}_t \doteq (d_t^1, d_t^2, \dots, d_t^N) \in \mathcal{D}$ to denote the current joint action of all agents in time slot t ; $p : \mathcal{Z} \times \mathcal{D} \times \mathcal{Z} \rightarrow [0, 1]$ is the state transition probability function, the joint action \vec{d}_t causes the transition from the game state z_t to the next new state z_{t+1} with a probability $p(z_{t+1}|z_t, \vec{d}_t)$, which should satisfy $\sum_{z_{t+1} \in \mathcal{Z}} p(z_{t+1}|z_t, \vec{d}_t) = 1$.

$\forall z_t \in \mathcal{Z}, \forall \vec{d}_t \in \mathcal{D}; r_i : \mathcal{Z} \times \mathcal{D} \rightarrow \mathbb{R}$ is the direct reward function of agent i , which evaluates the quality of state transition. Given the common state z_t , each agent i can perform an acceptable joint action $(d_t^1, \dots, d_t^i, \dots, d_t^N)$ and accumulate a reward $r_t^i(z_t, \vec{d}_t)$ immediately in new state z_{t+1} . Notably, the reward varies not only according to the current state and the action of agent i , but also the action choice of all other agents. Besides, state transitions obey the Markov property, where the sequence of events that precede the current state cannot determine the probability distribution of future states. An agent's reward and next state depend only upon the current state and the common joint decisions among all agents [39].

B. DMRE: Distributed MARL-Based Edge Caching Method

Since the distributed method is more effective in matching the edge caching system characteristics, in this part, Q -learning is expanded to the MG model, and DMRE is proposed to optimize the configuration of various caching strategies.

Indeed, the MG model can be seen as a set of matrix games associated with each state. In these games, each agent takes actions while conditioning on the behavior of others to satisfy the reward function, and $(\pi_1, \pi_2, \dots, \pi_N)$ is used to represent the joint strategy of N agents. Similar to Q -learning, the expected cumulative discounted reward of agent i is expressed by the state value function $V_i(z_t, \pi_1, \pi_2, \dots, \pi_N) = \mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r_t^i(z_t, \vec{d}_t)]$, where $\gamma \in (0, 1)$ is the discounting factor indicating the importance of the predicted future rewards. Each agent dedicates to learning a local π_i to maximize $V_i(z_t, \pi_1, \pi_2, \dots, \pi_N)$, meanwhile remaining robust to the actions of others over the remaining course of the game.

As the maximal state value $V_i(z_t, \pi_1^*, \dots, \pi_N^*)$ of an agent cannot be achieved by merely private strategy, and the return explicitly depends on the joint strategy of all agents, the concept of Nash equilibrium becomes crucial. Joint Nash equilibrium strategy $(\pi_1^*, \dots, \pi_i^*, \dots, \pi_N^*)$ can be considered

as an optimal solution of the multi-agent system, of which the game reaches the Nash equilibrium point and each agent's strategy π_i^* is the best response to the others. In other words, no agent can achieve a higher state value function by unilateral deviation, i.e., changing to any other strategy. Therefore, for $\forall z_t \in \mathcal{Z}, \forall i \in \mathcal{N}$, we formally have:

$$V_i(z_t, \pi_1^*, \dots, \pi_i^*, \dots, \pi_N^*) \geq V_i(z_t, \pi_1^*, \dots, \pi_i, \dots, \pi_N^*), \quad (5)$$

where $\forall \pi_i$ belongs to the set of strategies available to agent i .

Then, we desire a method to attain this equilibrium of the game without any prior knowledge about dynamics. Beforehand, we identify conditions that are more easily verifiable and extend the *Bellman Equation* to Nash equilibrium. Like Q -learning, we apply the dynamic programming principle to agent i 's reward $V_i(z_t, \pi_1^*, \dots, \pi_i^*, \dots, \pi_N^*)$ while the remaining policies $\pi_{\mathcal{N} \setminus \{i\}}^*$ are fixed, which will result in the temporal difference form as:

$$\begin{aligned} V_i(z_t, \pi_1^*, \dots, \pi_N^*) &= \max_{\pi_i} \left[r_t^i(z_t, \vec{d}_t) \right. \\ &\quad \left. + \gamma \sum_{z_{t+1} \in \mathcal{Z}} p(z_{t+1}|z_t, \vec{d}_t) V_i(z_{t+1}, \pi_i \cup \pi_{\mathcal{N} \setminus \{i\}}^*) \right]. \end{aligned} \quad (6)$$

Next, different from the single-agent Q -function² $Q_i(z_t, d_t)$, the Q -function in the multi-agent system varies to $Q_i(z_t, \vec{d}_t)$, which is interpreted to estimate the optimal expected sum of discounted rewards (state value function) for any individual agent i during the learning process, i.e., $V_i(z_t, \pi_1^*, \dots, \pi_N^*) = \max_{\pi_i} Q_i(z_t, \vec{d}_t)$. Specially, given the Nash equilibrium as a learning objective, $Q_i^*(z_t, \vec{d}_t)$ is referred to as a Nash Q -value for agent i when all agents follow a specified joint Nash equilibrium strategy from the next slot on. That is,

$$\begin{aligned} V_i(z_t, \pi_1^*, \dots, \pi_N^*) &= \text{Nash } Q_i(z_t, \vec{d}_t) \\ &= \pi_1^*(z_t, d_t^1) \cdots \pi_N^*(z_t, d_t^N) \cdot Q_i^*(z_t, \vec{d}_t), \end{aligned} \quad (7)$$

where $\pi_1^*(z_t, d_t^1) \cdots \pi_N^*(z_t, d_t^N) \cdot Q_i^*(z_t, \vec{d}_t)$ is a scalar, and

$$\begin{aligned} Q_i^*(z_t, \vec{d}_t) &= r_t^i(z_t, \vec{d}_t) + \gamma \sum_{z_{t+1} \in \mathcal{Z}} p(z_{t+1}|z_t, \vec{d}_t) V_i(z_{t+1}, \pi_1^*, \dots, \pi_N^*). \end{aligned} \quad (8)$$

As a result, payoffs for each agent i are equal to Q -value at this point, and the Nash equilibrium is rewritten in the following form:

$$Q_i^*(z_t) \pi_1^*, \dots, \pi_i^*, \dots, \pi_N^* \geq Q_i^*(z_t) \pi_1^*, \dots, \pi_i, \dots, \pi_N^*. \quad (9)$$

²For ease of presentation, we have not distinguished the concept of “state-action function” and “ Q -function” in this paper.

At time slot t , each agent i executes its joint action under the current state. After that, it observes its immediate reward, all other agents' actions and rewards, as well as the new state z_{t+1} . As for agents' actions choice and Q -values update, the Nash equilibrium reward is adopted here to replace the maximum reward iteration in the single-agent method. Therefore, we rely on the stagewise approach to learn strategies separately for every stage game (i.e., a time slot game, defined by interim Q -values). While inequality (9) is satisfied, the stage game $(Q_t^1(z_{t+1}), \dots, Q_t^i(z_{t+1}), \dots, Q_t^N(z_{t+1}))$ will be formed with the Nash equilibrium strategies $\pi_1^*(z_{t+1}), \dots, \pi_N^*(z_{t+1})$ under certain restrictions to Q -values' domain. To achieve this, each learning agent has to keep track and maintain a model of other agents' Q -values, and then update its own Q -values during the interaction. The strategies that constitute a Nash equilibrium can be interpreted as the optimal behavior under the current system state, while different methods for selecting among multiple Nash equilibria will, in general, yield different updates. Existing studies [40], [41] show that there always exists at least one equilibrium point in stationary strategies, and the learning protocol provably converges if every stage game has saddle points or an optimum global; agents are mandated to update in terms of these. Ultimately, each agent i chooses its actions by calculating a Nash equilibrium, and updates the Q -value with the following iterative formula:

$$\begin{aligned} Q_{t+1}^i(z_t, \vec{d}_t) &= (1 - \beta_t) Q_t^i(z_t, \vec{d}_t) + \beta_t [r_t^i(z_t, \vec{d}_t) + \gamma \text{Nash } Q_t^i(z_{t+1})], \end{aligned} \quad (10)$$

where $\text{Nash } Q_t^i(z_{t+1})$ is the Nash equilibrium of agent i under new state, $\beta_t \in (0, 1)$ is the learning rate parameter, indicating how far the current estimate $Q_t^i(z_t, \vec{d}_t)$ is adjusted toward the update target $r_t^i(z_t, \vec{d}_t) + \gamma \text{Nash } Q_t^i(z_{t+1})$. Each Q -function can definitely converge to the Nash Q -value through repeated interaction when an appropriate β is designed. The agent can derive the Nash equilibrium and choose its actions accordingly.

C. Complexity Analysis and Parameter Approximation for DMRE

According to the description above, each agent i can gradually learn the equilibrium at each stage, which is determined by the joint actions of all agents in the system. That means the agent needs to learn N Q -values of all agents and derive strategies from them. Specifically, agent i will update $(Q^1, \dots, Q^i, \dots, Q^N)$ for all system states and joint actions, while these Q -values are updated and maintained internally by it. Hence, DMRE is often computationally expensive for equilibrium computation, and there may not be enough memory to maintain the corresponding values in large scale scenario. Consequently, it is necessary to analyze the computation complexity and further give simplification of the proposed method.

Let $|\mathcal{Z}|$ denote the number of system states, and $|D_i|$ be the size of agent i 's action space D_i . Assuming $|D_1| = \dots = |D_N| = |D|$, the total number of entries in Q^i is $|\mathcal{Z}| \cdot |D|^N$.

As the agent has to maintain N Q -tables, the total space requirement is $N|\mathcal{Z}| \cdot |D|^N$. Then, the proposed method is linear in the number of states, polynomial in the number of actions, but exponential in the number of agents in terms of space complexity. On the other hand, the training time of DMRE is dominated by the calculation of the Nash equilibrium. That is, each agent will search for an optimal solution according to the current solutions of other agents while the computation process is nonlinear if the number of agents is more than two [39]. Thus, due to the remarkable efficiency and generalization, parameter approximation is introduced to reduce the computation complexity for DMRE in this paper, which legitimately simplifies the training targets.

Specifically, elements can be controlled by sequences approximately in the space that is composed of Q -functions. Correspondingly, we introduce a set of parameters θ and approximate any agent i 's Q -function by using the updated parameter. The Q -function is re-expressed as:

$$\hat{Q}^i(z, \vec{d}, \theta) \approx \theta^T \phi_i = \sum_{j=1}^N \theta_j \phi_{ij}(z), \quad (11)$$

where $\theta = (\theta_1, \theta_2, \dots, \theta_n)^T$ and the given eigenfunction $\phi = (\phi_1(s), \phi_2(s), \dots, \phi_n(s))^T$.

Since $\hat{Q}^i(z, \vec{d}) \doteq \hat{Q}^i(z, \vec{d}, \theta)$ is possible when the error is smaller, we attempt to find a set of parameters θ for which the corresponding approximate values $(\hat{Q}_\theta^1, \dots, \hat{Q}_\theta^N)$ approximately satisfy Eq. (8). Thus, the loss function is interpreted as the distance between the approximate value and true value of Q -functions, which can be converted to:

$$\begin{aligned} \text{Loss}(\theta) &= \mathbb{E}_{z \sim p(z_{t+1}|z_t, \vec{d}_t)} \left[(\hat{Q}^i(z, \vec{d}) - \hat{Q}^i(z, \vec{d}, \theta))^2 \right] \\ &= \frac{1}{T} \sum_{t=1}^T \left[r_t^i + \gamma \text{Nash } \hat{Q}_t^i(z_{t+1}, \theta) - \hat{Q}_t^i(z_t, \vec{d}_t, \theta) \right]^2. \end{aligned} \quad (12)$$

According to the *Bellman Equation* of the state-action functions, Stochastic Gradient Descent (SGD) is performed to train θ towards the target value by minimizing the loss function, i.e. $\text{argmin}(\text{Loss}(\theta))$. Therefore, the optimal value is generated by an error as:

$$\theta_{t+1} \approx \theta_t - \frac{1}{2} \beta_t \nabla \left(r_t^i + \gamma \text{Nash } \hat{Q}_t^i(z_{t+1}, \theta) - \hat{Q}_t^i(z_t, \vec{d}_t, \theta) \right)^2. \quad (13)$$

We then take a semi-gradient method and ignore the derivative of parameter θ with respect to the unknown Nash Q value. Thus, the update formula of parameter θ in state-action function $\hat{Q}^i(z, \vec{d}, \theta)$ of agent i is written as follows:

$$\begin{aligned} \theta_{t+1} &= \theta_t + \beta_t \left[r_t^i + \gamma \text{Nash } \hat{Q}_t^i(z_{t+1}) \right. \\ &\quad \left. - \hat{Q}_t^i(z_t, \vec{d}_t, \theta_t) \right] \nabla \hat{Q}_t^i(z_t, \vec{d}_t, \theta_t). \end{aligned} \quad (14)$$

Instead of updating the Q -function, parameter approximation can make it more efficient by using the updated

Algorithm 1 Parameter Approximation for DMRE

Input: agent set \mathcal{N} , state space \mathcal{Z} , joint action space \mathcal{D} , learning rate β , discount factor γ , exploration factor ϵ .

- 1 **Each agent** $i \in \mathcal{N}$ **do**
- 2 **Initialize** $\forall z \in \mathcal{Z}, \forall d_i \in D_1, \hat{Q}^i(z, \vec{d}, \theta) \leftarrow 0$, parameter θ .
- 3 **for** each episode **do**
- 4 Set $t = 0$, obtain the initial state z_0 by randomly caching.
- 5 **for** each slots of episode **do**
- 6 Derive an action d_t^i based on the ϵ -greedy strategy in the current z_t .
- 7 **for** $Q_i(z_t, \vec{d}_t) \neq Q_i^*(z_t, \vec{d}_t)$ **do**
- 8 Observe the rewards r_t^1, \dots, r_t^N , the joint actions \vec{d}_t , and the next state z_{t+1} ;
- 9 Update parameter θ according to Eq. (14);
- 10 Compute $Q^i(z_t, \vec{d}_t)$:
- 11
$$Q^i(z_t, \vec{d}_t, \theta) \leftarrow \theta^T \phi_i;$$
- 12 Solve Nash equilibrium strategies under state z_{t+1} ;
- 13 Let $t \leftarrow t + 1$.
- 14 Execute cache strategy according to π_i^* .

parameter θ . Throughout the training process of seeking the Nash equilibrium point, the agent first updates the target θ by minimizing the corresponding loss function (12) via SGD. Then, θ can continually train the Q -function towards the Nash equilibrium point. All of these processes make the agent approach the Nash equilibrium point effectively. More details of the simplified DMRE³ are summarized in **Algorithm 1**, and relevant theoretical proofs of convergence and feasibility are presented at the end of this paper.

D. DeepDMRE: Distributed Deep MARL-Based Edge Caching Method

DMRE is enabled to represent and update the parameter by creating a look-up table, which is essentially a tabular-based method, even though parameter approximation is used to simplify the training target. Parameter update rule Eq. (14) in an agent relies on the repeated calculation of Nash $\hat{Q}_t^i(z_{t+1})$ over all states, which generally performs inefficiently in large-scale scenarios. To circumvent the issue and make more expressive parametric models, we incorporate the technical advantage of Deep- Q Network⁴ into DMRE, and further develop a computationally efficient method (DeepDMRE) with neural network-based Nash equilibria approximation.

³As the simplified DMRE has lower computation complexity and is easier to implement than the original version, we only use the “DMRE” to denote the simplified version in the following.

⁴Deep- Q Network is a classical deep RL method based on value iteration, in which neural networks are adopted to estimate the Q -function.

As mentioned above, we can introduce a set of parameters and approximate any agent’s Q -function using Eq. (11). In contrast, DeepDMRE defines a specific model for the function approximation. For each parameter θ , we decompose $\hat{Q}^i(z, \vec{d}, \theta)$ into two components:

$$\hat{Q}(z, \vec{d}, \theta) = \hat{V}(z, \theta) + \hat{A}(z, \vec{d}, \theta), \quad (15)$$

where $\hat{V}(z, \theta)$ and $\hat{A}(z, \vec{d}, \theta)$ are produced by various neural network components. $\hat{V}(z, \theta)$ is the state value function under state z , and $\hat{A}(z, \vec{d}, \theta)$ is the advantage function of executing action under state z . The advantage function represents the optimality gap between \hat{Q} and \hat{V} . In particular, each agent i ’s behaviour must be the best response to the others when the action execution constitutes a Nash equilibrium. At this point, the advantage function $\hat{A}(z, \vec{d}, \theta) = \text{Nash } Q_i(z_t, \vec{d}_t) - V_i(z_t, \pi_1^*, \dots, \pi_N^*) = 0$.

In contrast to DMRE, DeepDMRE maintains neural network models with experience replay buffers on each agent. The basic concept underlying DeepDMRE is centralized training and distributed inference. During training sample collection, the agent will store its observed experience tuple $y_t = (z_t, \vec{d}_t, z_{t+1}, r_t)$ into the replay buffer, which involves its current state and available action set. The arrived experience samples are used to train the parameters of the function approximation in neural networks via SGD. Practically, this replaying enables the agent to randomly extract a minibatch of previous experience samples from the replay buffer for learning at each iteration. Empirical studies [15], [20], [33] have proved that the experience replay mechanism can enhance sample efficiency and break the correlation among data training.

Moreover, we generalize and adopt the model hypothesis proposed in [40], which considers a second order Taylor expansion $\hat{Q}^i(z, \vec{d}, \theta)$ in action around the Nash equilibrium. Together with the hypothesis, the advantage function for each agent i can be approximately written in a linear-quadratic form as:

$$\begin{aligned} & \hat{A}_i(z, \vec{d}, \theta) \\ &= \left(\vec{d}_{\mathcal{N} \setminus \{i\}} - \bar{\mu}_{\mathcal{N} \setminus \{i\}}^\theta \right)^\top \Psi_i^\theta(z) \\ & \quad - \left(\frac{d_i - \mu_i^\theta}{\vec{d}_{\mathcal{N} \setminus \{i\}} - \bar{\mu}_{\mathcal{N} \setminus \{i\}}^\theta} \right)^\top P_i^\theta(z) \left(\frac{d_i - \mu_i^\theta}{\vec{d}_{\mathcal{N} \setminus \{i\}} - \bar{\mu}_{\mathcal{N} \setminus \{i\}}^\theta} \right) \end{aligned} \quad (16)$$

where the block matrix $\Psi_i^\theta(x)$ and $P_i^\theta(x)$ are all deterministic probability distributions. Interested readers can find detailed information in [40].

This model hypothesis implicitly suggests that each agent’s Q -function can be approximately expressed as a linear-quadratic function of the actions. Each $\hat{Q}^i(z, \vec{d}, \theta)$ is a concave function of d_i , guaranteeing that Nash $\hat{Q}_t^i(z_{t+1})$ is bijective. Then, we can model $\hat{V}(z, \theta)$, $\bar{\mu}^\theta$, and relevant probability distributions via constructed neural networks, rather than modelling $\hat{Q}^i(z, \vec{d}, \theta)$. The Nash equilibrium is achieved when $\vec{d}^* = \bar{\mu}^\theta$, and the advantage function will become 0 simultaneously. Similar to the derivations in Eq. (7), we can simplify expressions of the value function and equilibrium

801 strategy as:

$$\hat{V}(z, \theta) = \text{Nash } \hat{Q}^i(z, \vec{d}, \theta) \quad (17)$$

802 and

$$\vec{\mu}(z) = \text{argmax } \text{Nash } \hat{Q}^i(z, \vec{d}, \theta) \quad (18)$$

803 Consequently, the loss function in Eq. (12) becomes more
804 tractable through this simplification. For a minibatch of J
805 sample, it can be converted as:

$$\begin{aligned} 806 \mathcal{L}(\theta) &= \frac{1}{J} \sum_{j=1}^J \mathcal{L}_j(\theta) \\ 807 &= \frac{1}{J} \sum_{j=1}^J \left(r_t^i + \gamma \hat{V}_t^i(z_{t+1}, \theta) - \hat{V}(z, \theta) - \hat{A}(z, \vec{d}, \theta) \right)^2. \end{aligned} \quad (19)$$

808 where each sample observation $y_j = (z_j, \vec{d}_j, z_{j+1}, r_j)$ is used
809 to train the parameter θ by minimizing the $\mathcal{L}_j(\theta)$.

810 Furthermore, instead of applying SGD with back-
811 propagation, we employ an actor-critic-based variant on the
812 loss function, enabling more stable and efficient convergence
813 processes. Specifically, there are N actor networks and one
814 centralized critic network, where each RSU adopts one of the
815 actor networks to execute its inference, and the MBS acts as a
816 centralized critic network to train the model and evaluate the
817 overall caching state. Since the value function can be modeled
818 independently through the decomposition in Eq. (15), an actor-
819 network is performed for parameters update by minimizing
820 the loss function in Eq. (19) over parameters θ . Notably, the
821 parameter θ gets separated as $\theta = (\theta_V, \theta_A)$ in this update rule,
822 where sub-parameter θ_V is used to model $\hat{V}(z, \theta_V)$, and sub-
823 parameter θ_A is used to model $\hat{A}(z, \vec{d}, \theta_A)$. Accordingly, the
824 training objective is to update these parameters by minimizing
825 the redefined loss function:

$$\mathcal{L}(\theta) = \frac{1}{J} \sum_{j=1}^J \mathcal{L}_j(\theta) = \frac{1}{J} \sum_{j=1}^J \hat{\mathcal{L}}(y_j, \theta_V, \theta_A), \quad (20)$$

826 where each agent minimizes $\mathcal{L}(\theta)$ by alternating between
827 minimization in sub-parameters θ_V and θ_A . Here, the loss
828 function of an individual sample $y_j = (z_j, \vec{d}_j, z_{j+1}, r_j)$
829 corresponds to the error in Eq. (8), which is expressed as:

$$\begin{aligned} 830 \hat{\mathcal{L}}(y_j, \theta_V, \theta_A) \\ 831 &= \left(r(z_j, \vec{d}_j) + \gamma \hat{V}(z_{j+1}, \theta_V) - \hat{V}(z_j, \theta_V) - \hat{A}(z_j, \vec{d}_j, \theta_A) \right)^2. \end{aligned} \quad (21)$$

832 More details of DeepDMRE are summarized in
833 **Algorithm 2**. With the locally linear-quadratic form of
834 the advantage function, we minimize the loss function in
835 Eq. (21) over sub-parameters θ_V and θ_A by actor-critic-based
836 iterative optimization and batch sampling. On the other hand,
837 we also apply and modify the ϵ -greedy based action
838 selection strategy for ensuring adequate exploration, where
839 ϵ is a decreasing parameter to achieve a tradeoff between
840 exploitation and exploration.

Algorithm 2 DeepDMRE: Distributed Deep MARL Based Edge Caching

841 **Input:** agent set \mathcal{N} , state space \mathcal{Z} , joint action space
842 \mathcal{D} , learning rate β , discount factor γ ,
843 exploration factor ϵ , experience replay buffer
844 \mathcal{M}_i , minibatch size J .

845 1 **Each agent** $i \in \mathcal{N}$ **do**

846 2 **Initialize** $\forall z \in \mathcal{Z}, \forall d_i \in D_i, \hat{Q}^i(z, \vec{d}, \theta) \leftarrow 0$, replay
847 buffer \mathcal{M}_i , neural network parameter (θ_V, θ_A) .

848 3 **for** each episode **do**

849 4 Set $t = 0$, obtain the initial state z_0 by randomly
850 caching.

851 5 **for** each slots of episode **do**

852 6 Derive an action $\vec{d} \leftarrow \vec{\mu}^\theta$ based on the ϵ -greedy
853 strategy in the current z_t .

854 7 **for** $Q_i(z_t, \vec{d}) \neq Q_i^*(z_t, \vec{d})$ **do**

855 8 Observe the rewards r_t^1, \dots, r_t^N , the
856 actions \vec{d}_t , and the next state z_{t+1} ;

857 9 Store observed experience tuple
858 $y_t = (z_t, \vec{d}_t, z_{t+1}, r_t)$ into the replay
859 buffer.

860 10 **if** $t \% \text{updateFrequency} == 0$ **then**

861 11 Randomly extract a minibatch of J
862 experiences $\{y_j\}_{j=1}^J$ from replay
863 buffer \mathcal{M}_i .

864 12 Update sub-parameters θ_V and θ_A by
865 minimizing $\mathcal{L}(\theta)$ as Eq. (20);

866 13 Perform a an actor-critic-based variant on
867 the loss function with respect to the
868 sub-parameters θ by $\frac{\partial \text{Loss}(\theta)}{\partial \theta}$;

869 14 Let $t \leftarrow t + 1$.

870 15 Update (θ_V, θ_A) .

VI. PERFORMANCE EVALUATION

871 This section conducts extensive simulations to assess the
872 performance of the proposed methods. Specifically, the su-
873 periority of the proposed methods is demonstrated over both
874 rule-based and learning-based benchmarks.

A. Simulation Setup

875 This paper considers an IoVs architecture composed of
876 one MBS and 5 RSUs, in which the coverage area of each
877 RSU is 200 m. RSUs are located at the center of each
878 region to serve the corresponding content requests, and the
879 initial cache capability of each RSU is limited to 100 MB.
880 Notably, the model can be easily extended to the case of
881 inconsistent RSU cache capacity. Besides, 10, 000 contents
882 are generated, and the size of each content is within the
883 range of [0.5, 1.5] MB. The content popularity follows an
884 MZipf distribution with $\theta = 0.73$. According to [4] and [35],
885 we set the channel gain as $30.6 + 36.7 \log_{10}(d)$, and the
886 noise power σ^2 as -95 dBm. The transmission delays between
887 MBS and RSU, and cooperative RSUs are 80 ms and 10 ms,

TABLE II
PARAMETERS SETUP

Parameter	Definition	Value
N	Default RSU number	5
T_2	Delay between any cooperative RSUs	10ms
T_3	Delay between the MBS and an RSU	80ms
P_i	Transmission power of RSUs	38dBm
B	Bandwidth of wireless links	10MHz
W	Bandwidth of optical fibre links	20MHz
σ^2	Power of background noise	-95 dBm
s_f	Size of requested contents	[0.5, 1.5] MB

865 respectively. Here, the price of the unit leased communication
866 resources (\$ / MB) is set as 0.005 for wireless links and
867 0.009 for optical links, while their bandwidths are 10 MHz
868 and 20 MHz, respectively. Meanwhile, we set the capacity of
869 experience replay buffer $\mathcal{M} = 500$, minibatch size $J = 32$,
870 and initial exploration factor $\epsilon = 0.05$. The learning rate
871 parameters for the actor-network and the critic-network are
872 set as 10^{-3} and 3×10^{-4} , respectively. All experiments are
873 run on 24 core - 2.4GHz Intel Xeon E5-2650 processor and
874 256GB RAM. The main parameters are listed in Table II.

875 For performance comparison, the following four benchmark
876 caching methods are introduced:

- 877 1) Least Frequently Used (LFU): The content with the least
878 requested times will be replaced firstly when an RSU's
879 cache capacity is full.
- 880 2) Least Recently Used (LRU): The content with the
881 longest unused time will be replaced firstly when an
882 RSU's cache capacity is full.
- 883 3) Q-learning [7]: An independent reinforcement learning
884 method, the process of cache replacement in each available
885 RSU is optimized individually.
- 886 4) Informed Upper Bound (IUB) [13]: An omniscient
887 method which assumes all of the requesting information
888 is perfectly known in advance. It provides an upper
889 bound of other methods and is usually hard to achieve
890 in practice.

891 Besides, the following metrics are used to evaluate the
892 performance of different methods quantitatively. 1) Content
893 access cost: the average total access cost of all requests in
894 each time slot. 2) Edge hit rate: the sum of requests served
895 by RSUs divided by the total number of requests. 3) Average
896 delay: the average access time for all requests in each time
897 slot.

898 B. Performance Comparison

899 1) *The convergence performance*: We plot the convergence
900 curve for three RL-based methods in Fig. 3, where the
901 number of contents and the cache capability of RSU are
902 fixed at 10,000 and 100MB, respectively. We can observe
903 that for these methods, each episode's total content access
904 cost decreases rapidly with the increase of the number of
905 episodes. The DeepDMRE converges to a relatively stable
906 value after running about 1800 episodes. Overall, it performs
907 best as it achieves a faster running time and obtains the same
908 convergence value as the original DMRE. On the other hand,

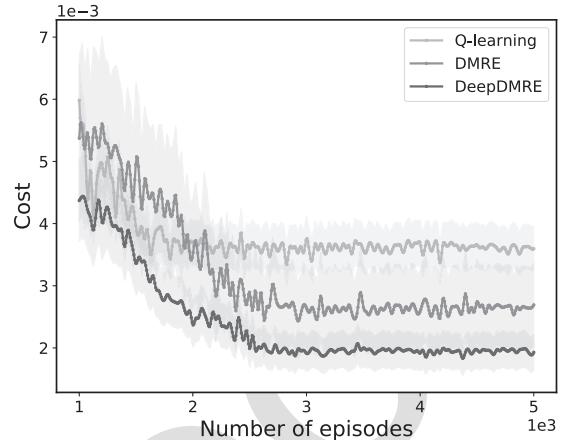


Fig. 3. Content access cost versus the number of episodes.

909 although the convergence speed of distributed MARL-based
910 methods (DMRE and DeepDMRE) are slightly slower than
911 that of the independent Q-learning based method, they can
912 reduce the content access cost up to 50% during the training
913 episodes.

914 2) *The Impact of Cache Capability of RSUs*: Then,
915 we explore the impact of different cache capabilities of RSUs
916 on the performance of different methods, as shown in Fig. 4.
917 The cache capacity of each RSU is changed from 100 MB
918 to 1000 MB and the number of contents is 10,000 in this case.
919 As expected, IUB performs best. However, it is impractical
920 since future information cannot be perfectly known. Among
921 the other caching methods, we observe that as the cache
922 capacity of RSUs increases, DeepDMRE consistently achieves
923 better performance with respect to the content access cost and
924 average delay. Specifically, DeepDMRE achieves the lowest
925 average delay of 33ms in the initial stage, and the content
926 access cost is reduced by about 47% and 26% compared with
927 independent Q-learning and DMRE, without mentioning the
928 simple rule-based methods LFU and LRU. In addition, for
929 all caching methods, the content access cost and the average
930 delay decrease with the increase of the cache capacity, which
931 is reasonable. The larger cache capacity enables the RSU to
932 cache more contents at the same time so that RSUs can meet
933 most content requests.

934 Similarly, the increase of cache capacity also has a positive
935 impact on edge hit rate. From Fig. 4(b), we can find that
936 the edge hit rate exhibits an increasing trend for all caching
937 methods when we increase the cache capacity of each RSU,
938 especially for DMRE and DeepDMRE, which are very close to
939 the value of IUB. This is because they considers the influence
940 of environments by other agents, and takes joint decisions to
941 learn the coordination and optimization of the whole system.
942 Notably, DMRE and DeepDMRE will generally sacrifice some
943 local hit rate and share cache capacity to serve requests from
944 neighboring RSUs. This tradeoff benefits much by reducing
945 the duplicate content transmission since MBS fetching is
946 largely avoided. Especially, DeepDMRE outperforms DMRE,
947 Q-learning, LFU, and LRU. For example, the edge hit rate is
948 improved by roughly 5%, 19%, 40%, and 35%, respectively,
949 when the cache capacity reaches 1,000 MB. At this moment,
950 content replacement processes may rarely occur in the RSUs.

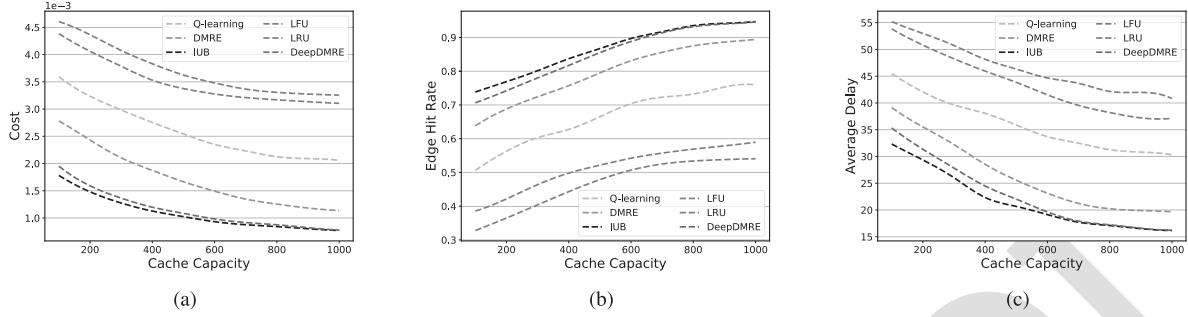


Fig. 4. (a) Content access cost versus the cache capacity of RSUs. (b) Edge hit rate versus the cache capacity of RSUs. (c) Average Delay versus the cache capacity of RSUs.

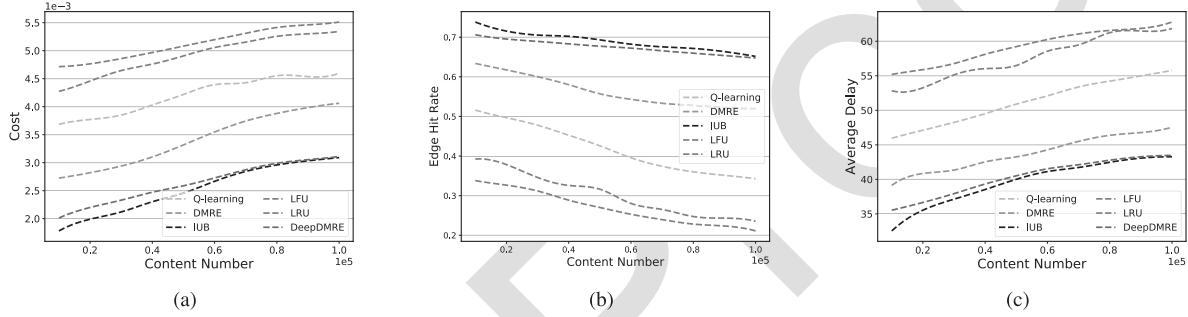


Fig. 5. (a) Content access cost versus the number of contents. (b) Edge hit rate versus the number of contents. (c) Average Delay versus the number of contents.

951 3) *The Impact of the Number of Contents in the System:*
952 This part changes the number of contents to evaluate the
953 performance of different methods. The number of contents is
954 varied from 10,000 to 100,000 and the initial cache capacity
955 of RSUs is set as 100 MB. As illustrated in Fig. 5, the content
956 access cost and the average delay of all methods increase
957 slightly as the number of contents increases. The increasing
958 trend does indicate that when the cache capacity is finite,
959 more contents will naturally cause frequent cache replacement,
960 as more contents need to be cached in RSUs now. As expected,
961 DeepDMRE performs better with a different numbers of
962 contents, except IUB. From Fig. 5(a) and (c), we can find
963 that even when there is a massive number of contents (i.e.,
964 100,000), DeepDMRE can still reduce the content access cost
965 by 11%, 16%, 28%, 34%, and the average delay by 15%, 21%,
966 33%, 37% compared to DMRE, Q-learning, LRU, and LFU,
967 respectively, which infers the effectiveness of DeepDMRE
968 under different numbers of contents.

969 Besides, the edge hit ratio of all methods decreases with
970 the increase of the number of contents. Among them, edge
971 hit rate of DeepDMRE is 64% when the number of contents
972 reaches 100,000. In contrast, LFU and LRU still perform
973 worst for this moment, and the edge hit rate is only 21%
974 and 23%, respectively. This situation may occur because LFU
975 and LRU learn only from one-step past and operate based
976 on simple rules, while RL-based edge caching methods can
977 be derived from the observed historical content demands and
978 concentrate more on the reward that agents can earn rather
979 than users' requests. Furthermore, the variation of edge hit
980 rate also confirms the analysis above. That is, DeepDMRE
981 improves the cache resources utilization significantly.

To summarize, quantitative results validate the effectiveness of DMRE and DeepDMRE under different scenarios. As observed, they not only reduce the content access cost, but also achieve a desirable edge hit rate and average delay.

VII. CONCLUSION

This paper first considered a cooperative edge caching supported IoVs architecture, which uses the cooperative caching between multiple RSUs and MBS to reduce the duplicate content transmission in the system. Then, the MDP was extended to the context of a multi-agent system, and the corresponding combinatorial multi-armed bandit problem was further tackled based on the framework of stochastic games. Specifically, a **Distributed MARL-based Edge** caching method (DMRE) was proposed firstly, where each agent can adaptively learn its best behavior in conjunction with other agents for intelligent caching. Meanwhile, we attempted to reduce the computation complexity of DMRE by parameter approximation, which legitimately simplifies the training targets. However, DMRE is enabled to represent and update the parameter by creating a lookup table, essentially a tabular-based method, which generally performs inefficiency in large-scale scenarios. To circumvent the issue and make more expressive parametric models, we combined the technical advantage of the Deep-Q Network into DMRE, and further developed a computationally efficient method, named DeepDMRE, by constructing neural networks for approximating the Nash equilibria. Simulation results demonstrated that DeepDMRE significantly outperforms other benchmark methods in different scenarios.

APPENDIX

This part gives the relevant theoretical proofs of convergence and the feasibility of the simplified DMRE as follows.

1015 A. Convergence Analysis

1016 Since the original DMRE has a provably convergence
 1017 performance with restrictive conditions [39], [41], the con-
 1018 vergence of the simplified DMRE can accordingly be proved
 1019 by demonstrating the effectiveness of SGD. Therefore, as θ
 1020 (which minimizes the corresponding loss function) is gradually
 1021 obtained via SGD, the state-action function (in the form of
 1022 $\hat{Q}^i(z, \vec{d}, \theta)$) will be updated towards the Nash equilibrium
 1023 point.

1024 For any state-action function $Q^k (k \in \mathcal{N})$ maintained
 1025 internally by a learning agent, we can define $f_k(\theta) =$
 1026 $(Q^k(z, \vec{d}) - \hat{Q}^k(z, \vec{d}, \theta))^2$. Then, the corresponding loss func-
 1027 tion is written as:

$$1028 L(\theta) \triangleq Loss(\theta) = \frac{1}{N} \sum_{k=1}^N f_k(\theta). \quad (22)$$

1029 As described in Section V-C, we train the parameter θ
 1030 towards the target value by minimizing the loss function
 1031 at each iteration. Then, owing to the numerical convex-
 1032 ity of the loss function, its gradient $\nabla L(\theta)$ satisfies the
 1033 *Lipschitz Continuity* with respect to the constant G , so that
 1034 for all θ_t and θ_{t+1} . We have:

$$1035 \|\nabla L(\theta_t) - \nabla L(\theta_{t+1})\| \leq G \|\theta_t - \theta_{t+1}\|. \quad (23)$$

1036 For a series of progressively decreasing learning rates β_t ,
 1037 the iterative formula of parameter θ can be expressed as:

$$1038 \theta_{t+1} = \theta_t - \beta_t \nabla L(\theta_t). \quad (24)$$

1039 Meanwhile, each θ in the training process should satisfy
 1040 $\max_k \|\nabla f_k(\theta_t)\| \leq B \|\nabla L(\theta_t)\|$, where B is a constant. That
 1041 is to say, the optimal θ for loss function is supposed to be a
 1042 stagnation point at the same time. We denote the error during
 1043 the iteration process as e_t , then the above Eq. (24) can be
 1044 rewritten with the following form:

$$1045 \theta_{t+1} = \theta_t - \beta_t (\nabla L(\theta_t) + e_t), \quad (25)$$

1046 where the error $e_t = \nabla f_k(\theta_t) - \nabla L(\theta_t)$.

1047 Additionally, since the mean of the error is equal to 0,
 1048 we therefore get:

$$1049 E[e_t] = E[\nabla f_k(\theta_t) - \nabla L(\theta_t)] = E[\nabla f_k(\theta_t)] - \nabla L(\theta_t) = 0. \quad (26)$$

1051 Furthermore, based on the derivations above, we can easily
 1052 get:

$$\begin{aligned} 1053 & E[\|e_t\|^2] \\ 1054 &= E[\|\nabla f_k(\theta_t) - \nabla L(\theta_t)\|^2] \\ 1055 &= E[\|\nabla f_k(\theta_t)\|^2 - 2 \langle \nabla f_k(\theta_t), \nabla L(\theta_t) \rangle + \|\nabla L(\theta_t)\|^2] \\ 1056 &= E[\|\nabla f_k(\theta_t)\|^2 - 2 \langle E[\nabla f_k(\theta_t)], \nabla L(\theta_t) \rangle + \|\nabla L(\theta_t)\|^2] \end{aligned}$$

$$\begin{aligned} 1057 &= \frac{1}{N} \sum_{k=1}^N [\|\nabla f_k(\theta_t)\|^2 - \|\nabla L(\theta_t)\|^2] \\ 1058 &\leq (B^2 - 1) \|\nabla L(\theta_t)\|^2. \end{aligned} \quad (27)$$

1059 Then, for inexact gradients, inequality (23) can be rewritten
 1060 by incorporating Eq. (24) with it:

$$1061 L(\theta_{t+1}) \leq L(\theta_t) + \langle \nabla L(\theta_t), \theta_{t+1} - \theta_t \rangle + \frac{G}{2} \|\theta_{t+1} - \theta_t\|^2. \quad (28)$$

1063 Also, we substitute error $e_t = \nabla f_k(\theta_t) - \nabla L(\theta_t)$ and
 1064 Eq. (24) into inequality (28), and further obtain:

$$\begin{aligned} 1065 L(\theta_{t+1}) &\leq L(\theta_t) - \beta_t \langle \nabla L(\theta_t), \nabla L(\theta_t) + e_t \rangle \\ 1066 &\quad + \frac{\beta_t^2 G}{2} \|\nabla L(\theta_t) + e_t\|^2 \\ 1067 &= L(\theta_t) - \beta_t \left(1 - \frac{\beta_t G}{2}\right) \|\nabla L(\theta_t)\|^2 \\ 1068 &\quad - \beta_t (1 - \beta_t G) \langle \nabla L(\theta_t), e_t \rangle + \frac{\beta_t^2 G}{2} \|e_t\|^2. \end{aligned} \quad (29)$$

1070 As the learning rate β_t tends to be infinitesimal, we take the
 1071 expectation of e_t on both sides of inequality (29). Then, a new
 1072 inequality (30), shown at the top of the next page, is derived
 1073 based on (26) and (27), as is shown at the bottom of this
 1074 paper. It shows that SGD can obtain the expected parameter
 1075 θ_t as the learning rate is small enough (satisfy $0 < \beta_t < \frac{2}{GB^2}$). At this point, we can get the corresponding state-action
 1076 function by the obtained parameter θ_t , the effectiveness of
 1077 SGD is therefore demonstrated. According to the previous
 1078 description, the convergence of the simplified DMRE is proved
 1079 theoretically.

1081 B. Feasibility Analysis

1082 To reduce the computation complexity, parameter approxi-
 1083 mation is introduced for DMRE in this paper, which legiti-
 1084 mately simplifies the training targets. However, a potential
 1085 error m may exist in the simplified DMRE, which is written
 1086 as:

$$1087 m = Q^i(z, \vec{d}) - \hat{Q}^i(z, \vec{d}, \theta) = Q^i(z, \vec{d}) - \sum_{j=1}^n \theta_j \phi_{ij}(z). \quad (31)$$

1088 According to inequality (30), the target parameter which
 1089 minimizes the loss function will be obtained via SGD through-
 1090 out the training process. Thus, the error possibly existing in
 1091 the simplified DMRE is contingent on the performance of
 1092 SGD actually. That is to say, the optimal parameter value can
 1093 minimize the corresponding error. On the other hand, an agent
 1094 needs to learn N Q -values of all agents and derive strategies
 1095 from them in the original DMRE, while these Q -values are
 1096 updated and maintained internally by it. This pattern makes
 1097 it obvious to get trouble in the curse of dimensionality with
 1098 the increase of agents. Moreover, the computation complexity
 1099 increases considerably, as well as the running time. To avoid
 1100 these bottlenecks, we further give an simplification of DMRE,
 1101 which legitimately transforms the training targets through
 1102 parameter approximation in Eq. (14).

$$\begin{aligned}
E[L(\theta_{t+1})] &\leq L(\theta_t) - \beta_t \left(1 - \frac{\beta_t G}{2}\right) \|\nabla L(\theta_t)\|^2 - \beta_t (1 - \beta_t G) \langle \nabla L(\theta_t), E[e_t] \rangle + \frac{\beta_t^2 G (B^2 - 1)}{2} \|\nabla L(\theta_t)\|^2 \\
&\leq L(\theta_t) - \beta_t \left(1 - \frac{\beta_t G}{2}\right) \|\nabla L(\theta_t)\|^2 + \frac{\beta_t^2 G (B^2 - 1)}{2} \|\nabla L(\theta_t)\|^2 \\
&\leq L(\theta_t) - \beta_t \left(1 - \frac{\beta_t G B^2}{2}\right) \|\nabla L(\theta_t)\|^2.
\end{aligned} \tag{30}$$

REFERENCES

- [1] P. Arthurs, L. Gillam, P. Krause, N. Wang, K. Halder, and A. Mouzakitis, "A taxonomy and survey of edge cloud computing for intelligent transportation systems and connected vehicles," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 7, pp. 6206–6221, Jul. 2022.
- [2] K. Jiang, C. Sun, H. Zhou, X. Li, M. Dong, and V. C. M. Leung, "Intelligence-empowered mobile edge computing: Framework, issues, implementation, and outlook," *IEEE Netw.*, vol. 35, no. 5, pp. 74–82, Sep. 2021.
- [3] Q. Luo, S. Hu, C. Li, G. Li, and W. Shi, "Resource scheduling in edge computing: A survey," *IEEE Commun. Surveys Tuts.*, vol. 23, no. 4, pp. 2131–2165, 4th Quart., 2021.
- [4] J. Zhang and K. B. Letaief, "Mobile edge intelligence and computing for the Internet of Vehicles," *Proc. IEEE*, vol. 108, no. 2, pp. 246–261, Feb. 2019.
- [5] H. Zhou, M. Li, N. Wang, G. Min, and J. Wu, "Accelerating deep learning inference via model parallelism and partial computation offloading," *IEEE Trans. Parallel Distrib. Syst.*, vol. 34, no. 2, pp. 475–488, Feb. 2023.
- [6] Y. Zhang, X. Lan, J. Ren, and L. Cai, "Efficient computing resource sharing for mobile edge-cloud computing networks," *IEEE/ACM Trans. Netw.*, vol. 28, no. 3, pp. 1227–1240, Jun. 2020.
- [7] K. Jiang, H. Zhou, D. Zeng, and J. Wu, "Multi-agent reinforcement learning for cooperative edge caching in Internet of Vehicles," in *Proc. IEEE 17th Int. Conf. Mobile Ad Hoc Sensor Syst. (MASS)*, Dec. 2020, pp. 455–463.
- [8] J. Cui, Y. Liu, Z. Ding, P. Fan, A. Nallanathan, and L. Hanzo, "Next-generation mm-Wave small-cell networks: Multiple access, caching, and resource management," *IEEE Veh. Technol. Mag.*, vol. 15, no. 1, pp. 46–53, Mar. 2020.
- [9] L. Su and V. K. N. Lau, "Data and channel-adaptive sensor scheduling for federated edge learning via over-the-air gradient aggregation," *IEEE Internet Things J.*, vol. 9, no. 3, pp. 1640–1654, Feb. 2022.
- [10] H. Li, K. Ota, and M. Dong, "Deep reinforcement scheduling for mobile crowdsensing in fog computing," *ACM Trans. Internet Technol.*, vol. 19, no. 2, pp. 1–18, May 2019.
- [11] H. Zhou, Z. Zhang, D. Li, and Z. Su, "Joint optimization of computing offloading and service caching in edge computing-based smart grid," *IEEE Trans. Cloud Comput.*, early access, Apr. 12, 2022, doi: 10.1109/TCC.2022.3163750.
- [12] W. Zhuang, Q. Ye, F. Lyu, N. Cheng, and J. Ren, "SDN/NFV-empowered future IoV with enhanced communication, computing, and caching," *Proc. IEEE*, vol. 108, no. 2, pp. 274–291, Feb. 2020.
- [13] W. Jiang, G. Feng, S. Qin, T. S. P. Yum, and G. Cao, "Multi-agent reinforcement learning for efficient content caching in mobile D2D networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 3, pp. 1610–1622, Mar. 2019.
- [14] H. Zhou, T. Wu, H. Zhang, and J. Wu, "Incentive-driven deep reinforcement learning for content caching and D2D offloading," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 8, pp. 2445–2460, Aug. 2021.
- [15] H. Zhang, M. Huang, H. Zhou, X. Wang, N. Wang, and K. Long, "Capacity maximization in RIS-UAV networks: A DDQN-based trajectory and phase shift optimization approach," *IEEE Trans. Wireless Commun.*, vol. 22, no. 4, pp. 2583–2591, Apr. 2023, doi: 10.1109/TWC.2022.3212830.
- [16] H. Zhang, N. Yang, W. Huangfu, K. Long, and V. C. M. Leung, "Power control based on deep reinforcement learning for spectrum sharing," *IEEE Trans. Wireless Commun.*, vol. 19, no. 6, pp. 4209–4219, Jun. 2020.
- [17] H. Zhou, Z. Zhang, Y. Wu, M. Dong, and V. C. M. Leung, "Energy efficient joint computation offloading and service caching for mobile edge computing: A deep reinforcement learning approach," *IEEE Trans. Green Commun. Netw.*, early access, Jul. 4, 2022, doi: 10.1109/TGCN.2022.3186403.
- [18] C. Zhong, M. C. Gursoy, and S. Velipasalar, "Deep multi-agent reinforcement learning based cooperative edge caching in wireless networks," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2019, pp. 1–6.
- [19] Y. Dai, D. Xu, K. Zhang, S. Maharjan, and Y. Zhang, "Deep reinforcement learning and permissioned blockchain for content caching in vehicular edge computing and networks," *IEEE Trans. Veh. Technol.*, vol. 69, no. 4, pp. 4312–4324, Apr. 2020.
- [20] N. Luong et al., "Applications of deep reinforcement learning in communications and networking: A survey," *IEEE Commun. Surveys Tuts.*, vol. 21, no. 4, pp. 3133–3174, 4th Quart., 2019.
- [21] X. Ye, M. Li, P. Si, R. Yang, Z. Wang, and Y. Zhang, "Collaborative and intelligent resource optimization for computing and caching in IoT with blockchain and MEC using A3C approach," *IEEE Trans. Veh. Technol.*, vol. 72, no. 2, pp. 1449–1463, Feb. 2023, doi: 10.1109/TVT.2022.3210570.
- [22] B. Jedari, G. Premsankar, G. Illahi, M. D. Francesco, A. Mehrabi, and A. Ylä-Jääski, "Video caching, analytics, and delivery at the wireless edge: A survey and future directions," *IEEE Commun. Surveys Tuts.*, vol. 23, no. 1, pp. 431–471, 1st Quart., 2021.
- [23] H. Wu and H. Lu, "Delay and power tradeoff with consideration of caching capabilities in dense wireless networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 10, pp. 5011–5025, Oct. 2019.
- [24] M. Yan, C. A. Chan, W. Li, L. Lei, A. F. Gygax, and I. Chih-Lin, "Assessing the energy consumption of proactive mobile edge caching in wireless networks," *IEEE Access*, vol. 7, pp. 104394–104404, 2019.
- [25] X. Zhang, Z. Qi, G. Min, W. Miao, Q. Fan, and Z. Ma, "Cooperative edge caching based on temporal convolutional networks," *IEEE Trans. Parallel Distrib. Syst.*, vol. 33, no. 9, pp. 2093–2105, Sep. 2022.
- [26] K. Zhang, J. Cao, S. Maharjan, and Y. Zhang, "Digital twin empowered content caching in social-aware vehicular edge networks," *IEEE Trans. Computat. Social Syst.*, vol. 9, no. 1, pp. 239–251, Feb. 2022.
- [27] J. Zhao, X. Sun, Q. Li, and X. Ma, "Edge caching and computation management for real-time Internet of Vehicles: An online and distributed approach," *IEEE Trans. Intell. Transp. Syst.*, vol. 22, no. 4, pp. 2183–2197, Apr. 2021.
- [28] W. C. Ao and K. Psounis, "Fast content delivery via distributed caching and small cell cooperation," *IEEE Trans. Mobile Comput.*, vol. 17, no. 5, pp. 1048–1061, May 2018.
- [29] Z. Gu, H. Lu, P. Hong, and Y. Zhang, "Reliability enhancement for VR delivery in mobile-edge empowered dual-connectivity sub-6 GHz and mmWave HetNets," *IEEE Trans. Wireless Commun.*, vol. 21, no. 4, pp. 2210–2226, Apr. 2022.
- [30] S. Zhang and J. Liu, "Optimal probabilistic caching in heterogeneous IoT networks," *IEEE Internet Things J.*, vol. 7, no. 4, pp. 3404–3414, Apr. 2020.
- [31] R. Wang, Z. Kan, Y. Cui, D. Wu, and Y. Zhen, "Cooperative caching strategy with content request prediction in Internet of Vehicles," *IEEE Internet Things J.*, vol. 8, no. 11, pp. 8964–8975, Jun. 2021.
- [32] 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, Jun. 2019.
- [33] 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, Nov. 2018.
- [34] G. Qiao, S. Leng, S. Maharjan, Y. Zhang, and N. Ansari, "Deep reinforcement learning for cooperative content caching in vehicular edge computing and networks," *IEEE Internet Things J.*, vol. 7, no. 1, pp. 247–257, Jan. 2020.

- [35] X. Wang, C. Wang, X. Li, V. C. M. Leung, and T. Taleb, "Federated deep reinforcement learning for Internet of Things with decentralized cooperative edge caching," *IEEE Internet Things J.*, vol. 7, no. 10, pp. 9441–9455, Oct. 2020.
- [36] H. Tian et al., "CoPace: Edge computation offloading and caching for self-driving with deep reinforcement learning," *IEEE Trans. Veh. Technol.*, vol. 70, no. 12, pp. 13281–13293, Dec. 2021.
- [37] X. Li, X. Wang, P.-J. Wan, Z. Han, and V. C. M. Leung, "Hierarchical edge caching in device-to-device aided mobile networks: Modeling, optimization, and design," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 8, pp. 1768–1785, Aug. 2018.
- [38] X. Chen, L. Jiao, W. Li, and X. Fu, "Efficient multi-user computation offloading for mobile-edge cloud computing," *IEEE/ACM Trans. Netw.*, vol. 24, no. 5, pp. 2795–2808, Oct. 2016.
- [39] J. Hu and M. P. Wellman, "Nash Q-learning for general-sum stochastic games," *J. Mach. Learn. Res.*, vol. 4, pp. 1039–1069, Nov. 2003.
- [40] P. Casgrain, B. Ning, and S. Jaimungal, "Deep Q-learning for Nash equilibria: Nash-DQN," 2019, *arXiv:1904.10554*.
- [41] Z. Zhang, D. Wang, and J. Gao, "Learning automata-based multiagent reinforcement learning for optimization of cooperative tasks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 10, pp. 4639–4652, Oct. 2021.



Shibo He (Senior Member, IEEE) is currently a Full Professor with Zhejiang University and also with the Alibaba-Zhejiang University Joint Research Institute of Frontier Technologies. His research interests include the Internet of Things, crowdsensing, and big data.

He served as the TPC Co-Chair for IEEE ScalCom 2014, the Finance and Registration Chair for ACM MobiHoc 2015, and the Symposium Co-Chair for IEEE ICC 2017 and IEEE GLOBECOM 2020. He serves on the editorial board for several journals, including IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY.

Huan Zhou (Member, IEEE) received the Ph.D. degree from the Department of Control Science and Engineering, Zhejiang University. He was a Visiting Scholar with Temple University from November 2012 to May 2013. He was a CSC supported Post-Doctoral Fellow with The University of British Columbia from November 2016 to November 2017. He is currently a Professor with the College of Computer, Northwestern Polytechnical University. He has published more than 70 research papers in some international journals and conferences, including IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS, IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTED SYSTEMS, IEEE TRANSACTIONS ON MOBILE COMPUTING, and IEEE TRANSACTIONS ON WIRELESS COMMUNICATIONS. His research interests include mobile social networks, VANETs, opportunistic mobile networks, and mobile data offloading. He was a TPC Member of IEEE MASS, ICCCN, GLOBECOM, ICC, and WCSP. He received the Best Paper Award of I-SPAN 2014 and I-SPAN 2018. He was the TPC Chair of EAI BDTA 2020 and the Local Arrangement Chair of I-SPAN 2018. He was a Lead Guest Editor of *Pervasive and Mobile Computing*. He is serving as an Associate Editor for EURASIP Journal on Wireless Communications and Networking and IEEE ACCESS.



Geyong Min (Senior Member, IEEE) received the B.Sc. degree in computer science from the Huazhong University of Science and Technology, China, in 1995, and the Ph.D. degree in computing science from the University of Glasgow, U.K., in 2003. He is currently a Professor of high performance computing and networking with the Department of Computer Science, College of Engineering, Mathematics, and Physical Sciences, University of Exeter, U.K. His research interests include computer networks, wireless communications, parallel and distributed computing, ubiquitous computing, multimedia systems, and modeling and performance engineering.

Kai Jiang (Student Member, IEEE) is currently pursuing the Ph.D. degree in cyberspace security with Wuhan University, China. His research interests include mobile edge computing, deep reinforcement learning, and the Internet of Vehicles.



Jie Wu (Fellow, IEEE) is the Chair and a Laura H. Carnell Professor with the Department of Computer and Information Sciences, Temple University, Philadelphia, PA, USA. Prior to joining Temple University, he was the Program Director of the National Science Foundation and a Distinguished Professor with Florida Atlantic University. He has regularly published in scholarly journals, conference proceedings, and books. His current research interests include mobile computing and wireless networks, routing protocols, cloud and green computing, network trust and security, and social network applications. He was a recipient of the 2011 China Computer Federation (CCF) Overseas Outstanding Achievement Award. He was the General Co-Chair/Chair of IEEE MASS 2006, IEEE IPDPS 2008, IEEE ICDCS 2013, ACM MobiHoc 2014, IEEE ICPP 2016, and IEEE CNS 2016, as well as the Program Co-Chair of IEEE INFOCOM 2011 and CCF CNCC 2013. He was the Chair of the IEEE Technical Committee on Distributed Processing (TCDP). He serves on several editorial boards, including IEEE TRANSACTIONS ON SERVICES COMPUTING and the *Journal of Parallel and Distributed Computing*. He was an IEEE Computer Society Distinguished Visitor and an ACM Distinguished Speaker. He is a CCF Distinguished Speaker.

1228	[35]	X. Wang, C. Wang, X. Li, V. C. M. Leung, and T. Taleb, "Federated deep reinforcement learning for Internet of Things with decentralized cooperative edge caching," <i>IEEE Internet Things J.</i> , vol. 7, no. 10, pp. 9441–9455, Oct. 2020.	1278
1229	[36]	H. Tian et al., "CoPace: Edge computation offloading and caching for self-driving with deep reinforcement learning," <i>IEEE Trans. Veh. Technol.</i> , vol. 70, no. 12, pp. 13281–13293, Dec. 2021.	1279
1230	[37]	X. Li, X. Wang, P.-J. Wan, Z. Han, and V. C. M. Leung, "Hierarchical edge caching in device-to-device aided mobile networks: Modeling, optimization, and design," <i>IEEE J. Sel. Areas Commun.</i> , vol. 36, no. 8, pp. 1768–1785, Aug. 2018.	1280
1231	[38]	X. Chen, L. Jiao, W. Li, and X. Fu, "Efficient multi-user computation offloading for mobile-edge cloud computing," <i>IEEE/ACM Trans. Netw.</i> , vol. 24, no. 5, pp. 2795–2808, Oct. 2016.	1281
1232	[39]	J. Hu and M. P. Wellman, "Nash Q-learning for general-sum stochastic games," <i>J. Mach. Learn. Res.</i> , vol. 4, pp. 1039–1069, Nov. 2003.	1282
1233	[40]	P. Casgrain, B. Ning, and S. Jaimungal, "Deep Q-learning for Nash equilibria: Nash-DQN," 2019, <i>arXiv:1904.10554</i> .	1283
1234	[41]	Z. Zhang, D. Wang, and J. Gao, "Learning automata-based multiagent reinforcement learning for optimization of cooperative tasks," <i>IEEE Trans. Neural Netw. Learn. Syst.</i> , vol. 32, no. 10, pp. 4639–4652, Oct. 2021.	1284
1235			1285
1236			1286
1237			1287
1238			1288
1239			1289
1240			1290
1241			1291
1242			1292
1243			1293
1244			1294
1245			1295
1246			1296
1247			1297
1248			1298
1249			1299
1250			1300
1251			1301
1252			1302
1253			1303
1254			1304
1255			1305
1256			1306
1257			1307
1258			1308
1259			1309
1260			1310
1261			1311
1262			1312
1263			1313
1264			1314
1265			1315
1266			1316
1267			1317
1268			1318
1269			1319
1270			1320
1271			1321
1272			1322
1273			1323
1274			1324