

Combining Lyapunov Optimization and Deep Reinforcement Learning for D2D Assisted Heterogeneous Collaborative Edge Caching

Ziyi Teng, *Student Member, IEEE*, Juan Fang[✉], *Member, IEEE*, Yaqi Liu, *Student Member, IEEE*,

Abstract—The problem of shared node selection and cache placement in wireless networks is challenging due to the difficulty of finding low-complexity optimal solutions. This paper proposes a new approach combining Lyapunov optimization and reinforcement learning (LoRL) to address content sharing in heterogeneous mobile edge computing (MEC) networks with base station (BS) and device-to-device (D2D) communication. Device in this network can choose to establish D2D links with neighboring devices for content sharing or send requests directly to the base station for content. Content access and energy consumption of shared nodes are modeled as a queuing system. The goal is to assign content sharing nodes to stabilize all queues while maximizing D2D sharing gain and minimizing latency, even in the presence of unknown network state distribution and user sharing costs. The proposed approach enables edge device to independently select associated nodes and make caching decisions, thereby minimizing time-averaged network costs and stabilizing the queuing system. Experimental results show that the proposed algorithm converges to the optimal policy and outperforms other policies in terms of total queue backlog trade-off and network cost.

Index Terms—Edge cache, content sharing, device-to-device communication, deep reinforcement learning, Lyapunov optimization.

I. INTRODUCTION

THE increasing number of smart devices joining wireless networks has led to a surge in wireless multimedia traffic [1]. However, a significant portion of this traffic consists of repeated requests for popular content such as news articles and TV shows. To address this issue, mobile edge computing (MEC) technology has emerged as a promising solution, allowing content retrieval from edge storage nodes like base stations or clouds. However, this approach often results in redundant data transmission within a short timeframe and is constrained by the cache capacity of edge devices.

To mitigate the limitations of storage capacity and duplication in edge devices, collaborative caching has been recognized as an effective solution [2]. The progression of integrated circuits has led to the integration of storage and computing abilities into edge devices. Consequently, content sharing via

device-to-device (D2D) communication has become feasible. Furthermore, the integration of D2D communication with MEC caching can yield additional benefits such as improved spatial frequency reuse and boosted cellular network throughput. These benefits can lead to reductions in transmission and backhaul loads while augmenting the service probability of tasks [3] [4].

Extensive research has focused on optimizing user performance in Mobile Edge Computing (MEC) networks through collaborative caching and device-to-device (D2D) caching. Content cache placement strategies have been specifically studied [5]–[11]. However, there are still unresolved issues that require attention. Firstly, existing studies mainly concentrate on improving cache performance through effective placement strategies. While this approach proves beneficial, there is a limitation on the amount of content a device can transfer within a given time frame. When simultaneous content requests exceed the transmission capacity, it causes delays and increases retrieval latency. To tackle this challenge, researchers propose using a virtual queuing system to represent access requests [10], [11]. This allows for optimization of queue management to mitigate transmission delays. Secondly, in collaborative networks, forwarding content sharing node requests incurs energy consumption costs. Virtual queues effectively represent node energy consumption dynamics. Therefore, ensuring stability in the consumption queue becomes crucial for overall system stability.

The selection of shared content delivery nodes and cached content replacement in a D2D-assisted MEC network are key issues. When a user's local cache cannot fulfill a request, it is necessary to determine which cache node (e.g., neighboring user node, local/base station, neighboring base station) should handle the request and how to cache the content [3]. User requests can be routed to other users or accessible edge nodes, and dynamic queues represent the content access and energy consumption of all edge nodes. The objective is to select shared nodes for unsatisfied requests, ensuring the stability of the request and energy consumption queues. However, there are three main considerations when deciding on shared nodes and caching. Firstly, the data rate for content sharing via D2D links depends on user distance and channel conditions. User preferences can be challenging to understand, but machine learning methods can be employed to learn and predict preferences for better decision-making. Secondly, content sharing incurs transmission costs, necessitating consideration of both

Manuscript received 27 Jun 2023; revised 27 October 2023; accepted 28 January 2024. This work is supported by Beijing Natural Science Foundation (4192007), and supported by the National Natural Science Foundation of China (61202076), along with other government sponsors.

Ziyi Teng, Juan Fang and Yaqi Liu are work with the Faculty of Information Technology, Beijing University of Technology, Beijing 100124, China (e-mail: tengziyi@emails.bjut.edu.cn; fangjuan@bjut.edu.cn; liuyq617@emails.bjut.edu.cn; Corresponding author: Juan Fang)

request delay and network cache node energy consumption to ensure network stability. Lastly, the cache replacement policy should adapt to the evolving wireless network environment, facilitating optimal replacement decisions based on changing network conditions.

Note that when confronted with a large number of users, simultaneously selecting shared nodes and making caching decisions for each user in the system poses a significant challenge. However, this challenge can be addressed by employing a stochastic Lyapunov optimization method that does not rely on prior knowledge of the network state distribution. Nevertheless, in our case, the uncertainty surrounding the user's mobile location, network channel state, and user preference further complicates the problem. This increased complexity arises from the need to stabilize all the queues while also considering cache decision problems. Therefore, We introduce context-aware preference learning strategies and propose dynamic shared node selection and cache replacement methods that combine Lyapunov optimization and reinforcement learning (LoRL). Thus, the main contributions of this paper are as follows:

- Our approach integrates Lyapunov optimization theory and reinforcement learning to develop a novel method for shared node selection and cache replacement. By considering random fading channels and data arrival, our method enables intelligent decision-making for user delivery node selection and cache replacement. The primary objective is to minimize user request latency and energy consumption while ensuring the stability of user request and energy consumption queues.
- We analyze the factors influencing the establishment of D2D links between users, considering both physical conditions and user preference similarity. To improve the effectiveness of D2D links, we propose a user preference model based on **AutoEncoder**. This model captures and measures the similarity in user preferences, enabling us to optimize the formation of D2D links based on shared interests.
- We integrate Lyapunov and reinforcement learning optimization techniques to address network optimization problems. Our approach involves deriving an upper bound for the drift plus cost function using Lyapunov optimization, which guides the training of the reinforcement learning model. By doing so, we enable the model to maximize user service quality while ensuring the stability of the system.

II. RELATED WORK

Existing research on edge cache optimization can be divided into: i) accuracy improvement of popularity prediction models; ii) joint optimization of edge caching and wireless resources; iii) collaborative caching.

A. Popularity Prediction

Due to the repetitive nature of content requests in the network, edge caching should cache content with high popularity. Cache placement policies based on popularity prediction have

demonstrated caching effectiveness, and reactive caching [12] [13] or proactive caching [14] [15] by analysing past historical request information to obtain request patterns has been extensively investigated. Hassine *et al* [16] used Auto-regressive and Moving Average (ARMA) models for centralised content popularity prediction. To overcome the sparse nature of user requests, chen *et al* [17] proposed a popularity prediction scheme based on weighted clustering and also described an explicit relationship between cache performance and popularity prediction accuracy. In addition, considering the private nature of user data, a federated learning approach is used for edge caching policy optimization [3] [18]. Due to the consumption caused by learning-based approaches, some online popularity prediction approaches that do not require a training phase have also been proposed [12] [19]. However, while the improvement of the accuracy of the popularity prediction model can improve the caching performance to a certain extent, the channel conditions as well as the network state in mobile edge networks can have an impact on the quality of service (QoS) of the users. Therefore, the relationship between cache performance and prediction accuracy is implicit, and the impact of popularity prediction errors on cache performance is difficult to estimate. The most popular (MP) algorithm with a priori popularity knowledge of the user request model is compared with the proposed algorithm in Qian *et al.* [21], and this is verified by the poor performance exhibited by the MP algorithm.

B. Joint Caching and Resource Optimization

The decentralization of cache capacity in MEC networks leads to strong coupling between cache strategies and wireless communication resource management. From the perspective of limited cache and wireless resources, study the caching problem. Existing research approaches to caching policies are classified as *optimization-based*, *reinforcement learning-based* and *deep learning-based*, and *game-theory-based*. **Optimization based** caching strategies are usually designed to maximize certain performance metrics within the constraints of network resources. For the optimization problem of complex joint wireless resources and caching, simple heuristic algorithms often require a long time and can only obtain sub-optimal solutions. Therefore, in existing research, a Lyapunov optimization method for online joint utility maximization and stability control framework has been proposed. This method decouples multi-stage stochastic optimization problems into continuous deterministic sub-problems for each stage, while providing theoretical guarantees for the long-term stability of the system [13] [20]. **Strategies based on reinforcement learning and deep learning** use observable user data or environmental states, such as user contextual information, channel gain or cache state, for online caching decisions and resource allocation. Wireless channels have a finite amount of data that can be transmitted per unit time, and proactive caching strategies are investigated in order to maximise bandwidth utilisation [21]–[24]. However, when the user request or environment state space is large, centralised Reinforcement Learning caching strategies are complex and difficult to handle, hence distributed reinforcement learning approaches are proposed [25]. **Game**

theoretic caching strategies have been used for caching and computing resource allocation in MEC network environments, where service providers or users compete among themselves for limited computing and bandwidth resources to meet their own interests [26], [27].

C. Collaborative Caching

In MEC networks, collaborative caching is an effective approach to reduce network service load, improve service latency, and enhance spectrum usage efficiency by expanding cache capacity. Existing research divides collaborative caching into two types based on cache location: *Coordinated Multi-Point (COMP)* and *D2D* caching. COMP involves obtaining requested content from adjacent devices, base stations, or other caching devices. The optimization goal in this context is to jointly optimize content caching and delivery decisions, considering network constraints and aiming to minimize service latency or content retrieval costs [28], [29]. To address the joint optimization problem of user collaboration nodes and cache placement, a decoupling approach can be employed for dual-scale joint optimization [30] [31].

Decentralised cooperative sharing methods address the challenges of diverse network states in wireless networks. They aim to solve the node selection problem in centralised cooperative transmission effectively [32]. These methods decentralise decision-making, allowing nodes to independently select cooperative partners and make transmission decisions based on local information. In networks with caching capabilities on the user side, collaborative sharing can be performed between D2Ds by establishing D2D communication [33]. Furthermore, the uncertainty of user requests and movement patterns makes it challenging to establish D2D connections. Therefore, based on learning methods, the user's movement trajectory as well as user request patterns are predicted to enable dynamic delivery of content [8]. Further, user data information is private and has the property of not being willing to be shared, while an effective local caching policy requires knowledge of user preference information. Therefore, to maximise the benefits for users, D2D content sharing approaches with social awareness and incentives have been proposed in existing studies [6] [27]. Moreover, in addition to horizontal collaboration between network cache nodes, vertical inter-tier collaboration between cache nodes is also an important solution to achieve service demand by expanding cache capacity. Similarly, some D2D-assisted heterogeneous collaboration approaches have been proposed to maximise spectrum efficiency and reduce request latency [3].

D. Our Contribution

Based on the aforementioned categorization, similar to the works in [3], [8], and [12], this paper investigates heterogeneous collaborative caching strategies supporting D2D assistance. Similar to the contributions in [14] and [15], we predict user preference popularity by analyzing historical request information. Similar to the contributions in [23], [24], [25], [26], [27], and [28], we utilize reinforcement learning for dynamic cache decision optimization. However, what sets our work

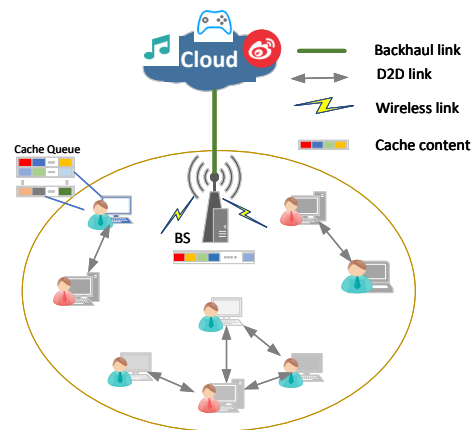


Fig. 1. D2D-assisted network architecture for heterogeneous collaborative edge caching.

apart from these contributions is that we employ predicted user preferences for D2D shared node selection. Furthermore, we combine Lyapunov optimization with reinforcement learning for user-associated node selection, cache decision-making, and maintaining stability in the request latency queue and cost consumption within the requesting network nodes.

III. SYSTEM MODEL

A. Network Model

Consider a wireless D2D-assisted heterogeneous collaborative network architecture, as shown in Fig.1. The network architecture includes three types of cache nodes, namely user equipment (UE), base stations (BS) and cloud servers. The cloud server connects to all BS via backhaul links to provide services to users, and the BS serve users via cellular links. In addition to the traditional BS-to-User use of the BS's wireless spectrum for content delivery, the considered network architecture allows for D2D links between users for content sharing. Given the abundant storage and computing resources of cloud servers, we assume that the cloud server has access to all the content that users may request within the storage area, denoted as $\mathcal{F} = \{1, 2, 3, \dots, F\}$, where F represents the total number of contents [2] [3] [5], and its size is denoted as $(s_f^{1 \times F})$. Each UE and BS has a restricted cache capacity to store content with high content popularity. The cache capacity of user device u is M_u , where $\forall u \in \mathcal{N} = \{1, 2, \dots, N\}$ is the set of tags of users. To simplify the model, we consider the existence of one BS in the network, serving the UE. In particular, the BS have the limited cache capacity, denoted as M_B .

To meet generality, the capacity of the network's cache nodes is $M_u \leq M_B \leq F$. All cache nodes in the network architecture, except for cloud nodes, are represented as $\mathcal{H} = \{\mathcal{N} \cup B\}$. In the network under consideration, operations are organized based on a time slot framework, represented as $\mathcal{T} = \{1, 2, 3, \dots, t\}$. The time axis is divided into equal time intervals, referred to as slots, with a small duration denoted as Δt and $t \in \mathcal{T}$. Within a time slot t , all network parameters (e.g., user location, channel quality,

content prevalence) remain constant. This time-slot-based organization enables the analysis and optimization of network performance within well-defined and consistent time intervals. Let us define the position of user u during time slot t as $l_{u,t}$, denoted by $l_{u,t} = \{x_{u,t}, y_{u,t}\}$. Here, $x_{u,t}$ and $y_{u,t}$ represent the coordinates of user u in the given time slot. It is essential to emphasize that within any time slot t , each user can make at most one request, and this request must be fulfilled prior to the commencement of the subsequent time slot.

B. D2D Sharing Mode

In the considered heterogeneous collaborative network, cache node content sharing has a restricted physical range, so it is assumed that the service range of the base station is R_B and the range of D2D is bounded by the radius R_u and satisfies $R_u \leq R_B$, i.e., $\pi R_u^2 \leq \pi R_B^2$, and cache nodes exceeding the service range cannot establish connections. In addition, user requests have the characteristics of being numerous and diverse. Therefore, in order to improve the effectiveness of D2D links, we will model the D2D connections from the physical domain and user similarity respectively in the following.

Physical domain: Due to physical limitations such as signal attenuation between D2D, only users within the user coverage can communicate [34]. Therefore, similar to [10], a graph $G_p = \{N, Y_p\}$ is introduced, where N represents the set of vertices of all users and $Y_p = \{(u, v) | e_{u,v}^p = 1, \forall u, v \in \mathcal{N}\}$ represents the edge set. When $e_{u,v}^p = 1$, it means that the distance between user equipment u and user equipment v is $L_{u,v}^t < R_u$, that is, within the communication range of the user, where $L_{u,v}^t = \sqrt{(x_{u,t} - x_{v,t})^2 + (y_{u,t} - y_{v,t})^2}$. Otherwise, $e_{u,v}^p = 0$.

User similarity: Humans are herd animals and tend to have a herd mentality for content acquisition, and each person has a different request preference P_u . Therefore, the concept of cosine similarity is introduced to represent the relationship between users' preferences. The cosine similarity between user u and user v is

$$S_{u,v}^v = \frac{P_u \cdot P_v}{\|P_u\| \|P_v\|}. \quad (1)$$

A user with a high preference similarity indicates that the content stored in the cache is more similar, and therefore, the content requested by the user is more likely to be stored.

User sharing probability: Based on the physical map and user similarity obtained above, we get the connection probability $R_{u,v} = e_{u,v}^p \cdot S_{u,v}^u$ between users, and define $R_u = (R_{u,v})_{1 \times N}$. Finally, we normalize R_u by $\frac{R_{u,v}}{\sum_{u \in \mathcal{N}} R_{u,v}} \rightarrow R_u$.

In the next sections, user preferences are predicted using an stack AutoEncoder (SAE)-based algorithm, and then content delivery and latency models for D2D-assisted heterogeneous networks are investigated in Section III-D.

C. User preference learning

In order to obtain the similarity between users, a hypothesis is made that, within a time frame, users send preference information to surrounding cache nodes. User preferences, i.e.,

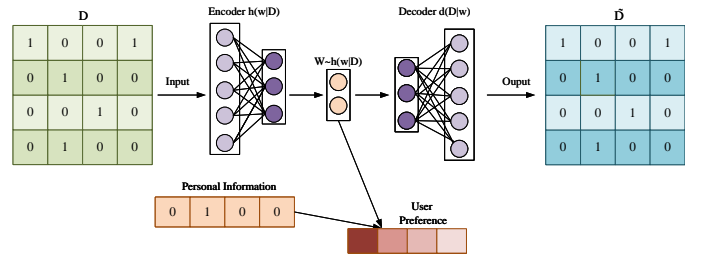


Fig. 2. A framework for learning user preferences based on contextual information.

the request popularity of a user, are uncertain and therefore obtaining numerical results for their popularity probability is very complex. Additionally, we assume that the preferences of users for the requested content follow an independent and identically distributed process. Unsupervised learning is a feasible approach to solve this problem. Therefore, in order to obtain the preference probability P_u of user u to further obtain the sharing probability $R_{u,v}$ between users, we propose a contextual information-based preference learning framework based on an unsupervised learning-hybrid filtering neural network model, as shown in Fig. 2.

The core of the proposed model is to minimise the discrepancy between the input and output, thus training a single hidden layer neural network to reconstruct the input data from the latent representations. The model consists of two main components: (i) an encoder that receives the input data and (ii) a decoder that outputs the results. The difference between the input and output is measured by the loss function F_L

$$F_L(D, \tilde{D}) = \sum_{i=1}^d (\tilde{D}_i - D_i), \quad (2)$$

where $\{D_1, D_2, \dots\}$ is the input dataset, $D_{(i)} \in \mathbb{R}^d$ represents the dimension d of each element in the dataset, mainly including the user's contextual information, such as name, gender, age, movie, movie rating and etc. The input data is represented implicitly using the encoder through the activation function $h(x)$ for the input mapping as $h_{w,b}(D_i) \approx \tilde{D}_i$, where $\{\tilde{D}_1, \tilde{D}_2, \dots\}$ is the output of the corresponding encoder and the parameters w and b are the weight matrix and bias vector, respectively.

The training of user preferences means that the model is to be continuously updated and used to minimise the reconstruction error of the input data set. By training the SAE neural network, the hidden features encoding z of the training data are obtained and these features are used to calculate the user similarity. Since what content is needed in the future, i.e., the popularity of the user's request, may depend on the user's context. Therefore, by combining the hidden feature encoding z with the user's contextual information, the user's preferences, i.e., the user's own content request popularity, are obtained.

D. Content Transmission and Delay Model

In a heterogeneous collaborative network, a user's cache list in time slot t can be represented as $x_u = \{x_{u,f} \in \{0, 1\}, u \in$

$\mathcal{N}, f \in \mathcal{F}$. $x_{u,f} = 1$ means that content f is stored in the cache, otherwise, if $x_{u,f} = 0$, it means that the user does not store content f . When the user's requested content f cannot be satisfied locally, the user selects a node for content sharing from the network. We define $b_{u,t} = \{b_{u,t}^n \in \{0, 1\} | u \in \mathcal{N}, n \in \mathcal{H}/u\}$ to denote the set of associated candidate nodes for user u . $b_{u,t}^n = 1$ indicates that user u selects cache node n to process user requests. Specifically, $b_{u,t}^0 = 1$ indicates that users directly obtain content requests from the base station. Therefore, in D2D assisted heterogeneous networks, there are several methods for obtaining content.

- 1) The user's request is saved in the local cache list, i.e., the request is satisfied in the local cache. Therefore, the corresponding user request delay $d_{u,t}^L$ is 0.
- 2) If a user request is not saved in the local cache list, we can establish a D2D link to obtain the request from nearby users n . In line with the approach taken in many existing studies [3][5][6], we adopt orthogonal models to allocate non-overlapping radio resources for D2D transmission. This allocation scheme involves dividing the bandwidth of each node into equal sub-bands and assigning equal sub-bands to each node while ensuring no interference among them. Thus, the request delay $d_{u,t}^{D2D} = \frac{s_f}{r_{u,t}^n}$, where s_f denotes the size of the user request file and $r_{u,t}^n$ denotes the data transfer rate between user u and user n . $r_{u,t}^n = B^D \log_2(1 + g_t^D h_t^D / L_{u,t}^n)$, with parameters B^D , g_t^D , h_t^D denote the inter-user channel bandwidth, transmission power consumption and channel gain, respectively.
- 3) The user u can also send a request to the local base station for content retrieval. If the base station can fulfill the request, the delay for file f is given by $d_{u,t}^B = \frac{s_f}{r_{u,t}^B}$. The transmission rate is calculated using the formula $r_{u,t}^B = B^B \log_2(1 + \frac{g_t^B h_t^B}{L_{u,t}^B})$, where $L_{u,t}^B$ represents the path loss between user u and the base station at time slot t . Here, $r_{u,t}^B$ represents the transmission rate from user u to the base station at time slot t . It is determined by parameters such as the bandwidth B^B , transmission power consumption g_t^B , and channel gain h_t^B associated with the user-to-base station link.
- 4) Eventually, when content cannot be obtained through direct sharing or local caching, the base station forwards the request to the cloud server. The latency of fetching the content from the cloud server is the sum of the base station delay $d_{u,t}^B$ and the transmission time required to transfer the content from the cloud server. This can be expressed as $d_{u,t}^c = d_{u,t}^B + \frac{s_f}{r_c}$, where s_f is the size of the requested content and r_c is the constant transmission rate between the base station and the cloud server.

Therefore, the request latency for the user of the above fetch request method satisfies $d_{u,t}^L \leq d_{u,t}^{D2D} \leq d_{u,t}^B \leq d_{u,t}^c$. When the user makes the selection decision $b_{u,t}^n$, if the selected content sharing cache node cannot meet the request, it will obtain the request from the same or upper layer cache node by the above collaborative method.

IV. QUEUE MODEL AND CACHING PROBLEM

A. Queuing Model

When user u initiates a request, if the request cannot be fulfilled locally, a content sharing node $b_{u,t}^n$ is selected. Based on the chosen delivery node $b_{u,t}^n$, the content access queue associated with node n is denoted as $Q_n(t)$. The dynamic backlog of all queues in the network is captured by the vector $Q(t) = \{Q_n(t) | n \in \mathcal{H}\}$. The arrival rate of $Q_n(t)$ represents the total size of content shared by node n as selected by users. Therefore, the arrival rate of queue is

$$A_n(b_{u,t}^n) = \sum_{u \in \mathcal{N}, n \in \{\mathcal{H}/u\}} b_{u,t}^n s_f. \quad (3)$$

The service rate of the queue $Q_n(t)$ is expressed as the average rate at which the current request is satisfied. Based on the above description, the dynamic request queue $\{Q_n(t)\}_{n=1}^{\mathcal{H}}$ at any time slot t can be dynamically described as

$$Q_n(t+1) = \max\{Q_n(t) + A_n(b_{u,t}^n) - r_{u,t}^n, 0\} \quad (4)$$

Specifically, at the initial stage, the request queue maintained by user n is set to $Q_n(0) = 0$. In addition, since the system state is random, the system-dependent queuing vector process $\{Q_n(t)\}_{t \in \mathcal{T}}$ is also random.

When users select nodes, there is a cost associated with content sharing. This cost is related to the delivery rate of the content from the user to the selected node. We assume that the transmission cost is proportional to the data rate, which is a general performance metric that can be converted to other metrics such as battery life, transmission delay, and interference. Following similar calculations in prior studies [10], we assume that the transmission cost is charged per unit of data rate. Therefore, at any given time slot t , the transmission cost from the selected node n to the user u can be expressed as:

$$C_n(b_{u,t}^n) = b_{u,t}^n \cdot p_{u,n}^{trans} \cdot r_{u,t}^n, \quad (5)$$

where $p_{u,n}^{trans}$ is denoted as the transmission power consumption from user u to node n . To deal with the network cost of the delivery nodes in the network, virtual cost queues $Y(t) = \{Y_n(t) | n \in \mathcal{H}\}$ are introduced. Specifically, at the initial stage, the cost queue maintained by user n is set to $Y_n(0) = 0$. Based on the above description, at any time slot t , the dynamic cost queue $\{Y_n(t)\}_{n=1}^{\mathcal{H}}$ can be dynamically described as

$$Y_n(t+1) = \max\{Y_n(t) + C_n(b_{u,t}^n) - \epsilon v, 0\}, \quad (6)$$

where ϵ is a positive scaling factor. Dynamic queue $Y_n(t)$ can be seen as a random energy consumption $C_n(b_{u,t}^n)$ and a fixed service rate ϵv .

B. Problem Formulation

In this paper, the overall objective is to make dynamic node selection and caching decisions that maximise D2D sharing gains while minimising user request latency. The specific problem under consideration can be formulated as

- 1) **D2D shared gain:** If the user's request cannot be satisfied locally (i.e., $x_{u,f} = 0$), then by selecting the content sharing node $b_{u,t}^n$ establishes a D2D link with adjacent user n to meet user requests, where $n > 0$. Therefore, the gain obtained through the D2D sharing method can be expressed as

$$G_{u,1}^t = \sum_{n \in \mathcal{N}} s_f \cdot R_{u,n} \cdot b_{u,t}^n \cdot x_{n,f} \cdot (1 - x_{u,f}), \quad (7)$$

$$\forall u \in \mathcal{N}, \forall n \in \mathcal{H}, \forall t \in \mathcal{T}, \forall f \in \mathcal{F}.$$

- 2) **Content fetch gain:** To account for the fact that a User Equipment (UE) can only share content with a chosen node during a time slot, we introduce an average queue delay that is directly proportional to the serving UE. Consequently, the network's benefit from user acquisition requests can be expressed as follows:

$$G_{u,2}^t = \begin{cases} \tau e^{-d_{u,t}^L}, & \text{Local Cache} \\ \tau e^{-d_{u,t}^{D2D}}, & \text{D2D Communication} \\ \tau e^{-d_{u,t}^B}, & \text{Communication to BS} \\ \tau e^{-d_{u,t}^C}, & \text{Cloud Service} \end{cases} \quad (8)$$

Where τ represents the introduced parameter, we observe a negative exponential function that highlights the inverse relationship between user request delay and channel gain. In simpler terms, when the user request delay is low, a larger gain is obtained. By using this formulation, we can analyze the impact of user request delay on the network's benefit.

Based on the above description, the gain of the obtained content is represented as

$$\mathbb{G}(b_{u,t}^n, x_{u,f}) = \sum_u^N \{\lambda_1 G_{u,1}^t + \lambda_2 G_{u,2}^t\}, \quad (9)$$

where λ_1 and λ_2 is the two introduced parameters, which satisfy the condition $\lambda_1 + \lambda_2 = 1$, where $0 \leq \lambda_1, \lambda_2 \leq 1$. These parameters represent the weights or proportions assigned to D2D gain and delay, respectively. Therefore, the goal of this article is to optimize the problem \mathcal{P}_1 , can be expressed as

$$\max \mathbb{G}(b_{u,t}^n, x_{u,f}) \quad (10)$$

For ease of understanding, the notation used in this article is summarized in Table I.

V. DYNAMIC CONTENT CACHE AND NODE SELECTION ALGORITHM

A. Stochastic Lyapunov Optimization

Directly solving the aforementioned problem \mathcal{P}_1 becomes a challenging task without prior knowledge of the system's status, queue backlog, and network cost distribution. In addition, the choice of content sharing delivery nodes creates an imbalance in the request queue and delivery cost of caching nodes in the network. In order to tackle the aforementioned challenges, we employ stochastic Lyapunov optimization methods

TABLE I
MODELING PARAMETERS AND NOTATIONS

| | |
|--|--|
| t | Index of the time slot. |
| $\mathcal{N} = \{1, 2, \dots, N\}$ | Set of user labels. |
| $\mathcal{H} = \{\mathcal{N} \cup \mathcal{B}\}$ | All the cache node set in network and B is base station node. |
| $\mathcal{F} = \{1, 2, \dots, F\}$ | Set of all contents. |
| M_i, M_B | Cache capacities of UE- i and the MEC server, respectively. |
| $l_{u,t} = \{x_{u,t}, y_{u,t}\}$ | User u position in time slot t and $x_{u,t}$ and $y_{u,t}$ is the coordinate of user u , respectively. |
| R_u, R_B | Service scope of users and base stations. |
| P_u | User preferences. |
| g_i^n, g_t^B | Channel gain. |
| S_u^v | Preference similarity between user u and user v . |
| x_u | User cache state. |
| $b_{u,t}^n = 1$ | Indicates that user u selects cache node n to process user requests. |
| s_f | Request the size of file f . |
| $r_{u,t}^n$ | The transmission rate processed by node n for user u . |
| $d_{u,t}^L, d_{u,t}^{D2D}, d_{u,t}^B, d_{u,t}^C$ | Transmission delay. |
| $A_n(b_{u,t}^n)$ | Arrival rate of queue n . |
| $Q(t), Y(t)$ | Dynamic of the queue backlog and virtual cost queue for delivery nodes. |
| $h_u(t)$ | Channel gain between user u and other nodes. |
| $C_n(b_{u,t}^n)$ | The transmission cost from the selected node n to the user u . |
| $p_{u,n}^{trans}$ | Transmission power consumption from user u to node n . |
| ϵv | fixed service rate. |

to regulate the selection of content sharing nodes. The primary objective of this approach is to minimize the average network latency while ensuring the stability of both the request and cost queues within the network, all without relying on any prior knowledge. We started by defining the functions used in our analysis.

Definition 1 (Quadratic Lyapunov Function): In order to jointly control the request and network cost queues of any time slot t , the total queue is defined as $Z(t) = \{Q_n(t), Y_n(t)\}$, where $\{Q_n(t)\}_{n=1}^{\mathcal{H}}$ and $\{Y_n(t)\}_{n=1}^{\mathcal{H}}$. The quadratic Lyapunov function $L(Z(t))$ of the random queuing process is equal to half the sum of the squares of the backlogs of all current queues, which is

$$L(Z(t)) = \frac{1}{2} \left(\sum_{n \in \mathcal{H}} Q_n(t)^2 + \sum_{n \in \mathcal{H}} Y_n(t)^2 \right). \quad (11)$$

The Lyapunov function is a scalar measure of the total queue backlog in the network, with a smaller $L(Z(t))$ indicating a lower queue occupancy in the network.

Definition 2 (Conditional Expectation Lyapunov Drift): At any time slot t , the conditional expectation Lyapunov drift ΔH_t represents the expectation of the time slot difference of the Lyapunov function, namely

$$\Delta H_t = E\{L(Z(t+1)) - L(Z(t)) | Z(t)\}, \quad (12)$$

which ΔH_t describes the variation of the quadratic Lyapunov function, i.e., the degree of fluctuation of the function. A smaller ΔH_t indicates a more stable queue in the network system. Therefore, we choose to minimise ΔH_t for each time slot t to stabilise the whole system network. However, if one wishes to stabilize the request and energy consumption

queues in the network while minimizing average latency, one must add the expected cost $E\{\mathbb{G}(b_{u,t}^n, x_{u,f})\}$ to ΔH_t [10] and then, transform the function that minimises ΔH_t , into a cost function that maximises the drift-plus-cost

$$V \cdot E\{\mathbb{G}(b_{u,t}^n, x_{u,f})|Z(t)\} - \Delta H_t, \quad (13)$$

where the weights $V \geq 0$ to balance the impact on network cost and network stability. In the following, to obtain an upper bound on the drift plus cost in (13) within an arbitrary time slot t . Firstly, we have

$$Q_n(t+1)^2 = Q_n(t)^2 + 2Q_n(t)(A_n(b_{u,t}^n) - r_{u,t}^n) + (A_n(b_{u,t}^n) - r_{u,t}^n)^2, \quad (14)$$

$$Y_n(t+1)^2 = Y_n(t)^2 + 2Y_n(t)(C_n(b_{u,t}^n) - \epsilon v) + (C_n(b_{u,t}^n) - \epsilon v)^2. \quad (15)$$

Secondly, by combining equation (4)(6), equation (12) can be rewritten as

$$\begin{aligned} \Delta H_t = & \sum_{n \in \mathcal{H}} \frac{1}{2} E\{(A_n(b_{u,t}^n) - r_{u,t}^n)^2 + Q_n(t) \cdot \\ & ((A_n(b_{u,t}^n) - r_{u,t}^n)|Q_n(t)) + \frac{1}{2}(C_n(b_{u,t}^n) - \epsilon v)^2 \\ & + Y_n(t) \cdot ((C_n(b_{u,t}^n) - \epsilon v)|Y_n(t))\}. \end{aligned} \quad (16)$$

Therefore, at any slot t , the drift-plus-cost function in (13) is upper-bounded by

$$\begin{aligned} V \cdot E\{\mathbb{G}(b_{u,t}^n, x_{u,f})|Z(t)\} - \Delta H_t \leq & V \cdot E\{\mathbb{G}(b_{u,t}^n, x_{u,f})|Z(t)\} \\ & - (B + \sum_{n \in \mathcal{H}} Y_n(t) \cdot ((C_n(b_{u,t}^n) - \epsilon v)|Y_n(t)) + \sum_{n \in \mathcal{H}} Q_n(t) \cdot \\ & ((A_n(b_{u,t}^n) - r_{u,t}^n)|Q_n(t))), \end{aligned} \quad (17)$$

where B is a constant independent of V and the $Q_n(t)$ and $Y_n(t)$ in the total queue $Z(t) = \{Q_n(t), Y_n(t)\}$ are independent of each other. Therefore, define $B = B_1 + B_2$, where B_1 and B_2 can be obtained separately from

$$\begin{aligned} \frac{1}{2} E\{\sum_{n \in \mathcal{H}} (A_n(b_{u,t}^n) - r_{u,t}^n)^2\} & \leq \frac{1}{2} \sum_{n \in \mathcal{H}} E[(A_n(b_{u,t}^n))^2 + (r_{u,t}^n)^2] \\ & \leq \frac{1}{2} N \cdot (s_f)^2 + \sum_{n \in \mathcal{H}} (r_{u,t}^n)^2 = B_1, \end{aligned} \quad (18)$$

$$\begin{aligned} \frac{1}{2} E\{\sum_{n \in \mathcal{H}} (C_n(b_{u,t}^n) - \epsilon v)^2\} & \leq \frac{1}{2} \sum_{n \in \mathcal{H}} E[(C_n(b_{u,t}^n))^2 + (\epsilon v)^2] \\ & \leq \frac{1}{2} \sum_{n \in \mathcal{H}} (p_{u,n}^{trans} \cdot r_{u,t}^n)^2 + N \cdot (\epsilon v)^2 = B_2. \end{aligned} \quad (19)$$

The terms of the second inequality in (18) and (19) relate to the content sharing node $b_{u,t}^n$ selected by user u . Instead of minimizing the drift-plus-cost function in (13), we minimize its upper-bound function. Therefore, in order to minimize the

right-hand side of the inequality (17), it is necessary to consider the current historical request queue $Q(t)$ and energy cost queue $Y(t)$ in time slot t . This can be achieved by selecting the appropriate shared delivery Node b_i^n . Then, we obtain the following optimization problem $\mathcal{P}_2 \max \Omega(b_{u,t}^n, x_{u,f}|Z(t))$, where the objective function

$$\begin{aligned} \Omega(b_{u,t}^n, x_{u,f}|Z(t)) = & V \cdot E\{\mathbb{G}(s^t, a^t)|Z(t)\} - \sum_u \sum_{n \in \mathcal{H}}^N \\ & Q_n(t) \cdot A_n(b_{u,t}^n|s_u^t) - Y_n(t) C_n(b_{u,t}^n|s_u^t) = \\ & \sum_u^N (\sum_{n \in \mathcal{H}} (s_f \cdot R_{u,n} \cdot b_{u,t}^n \cdot x_{n,f} \cdot (1 - x_{u,f}) + \tau e^{-d_{u,t}^n} - \\ & Q_n(t) \cdot b_{u,t}^n \cdot s_f - Y_n(t) \cdot b_{u,t}^n \cdot p_{u,n}^{trans} \cdot r_{u,t}^n)). \end{aligned} \quad (20)$$

$$s.t. \quad b_{u,t}^n, n \in \{\mathcal{N} \cup B\}/u, u \in \mathcal{N} \quad (19a)$$

$$\sum_{n \in \{\mathcal{N} \cup B\}} b_{u,t}^n = 1 \quad (19b)$$

B. Deep Reinforcement Learning for Shared Delivery Node Selection and Cache Replacement

Notice that the problem \mathcal{P}_2 is linear and, therefore, decomposable. In particular, we can then decompose this problem into N subproblems, given by

$$\begin{aligned} \Omega^u(b_{u,t}^n, x_{u,f}|Z(t)) = & \sum_{n \in \mathcal{H}} (V \cdot (s_f \cdot R_{u,n} \cdot b_{u,t}^n x_{n,f} \cdot (1 - x_{u,f}) \\ & + \tau e^{-d_{u,t}^n}) - Q_n(t) \cdot b_{u,t}^n \cdot s_f - Y_n(t) \cdot b_{u,t}^n \cdot p_{u,n}^{trans} \cdot r_{u,t}^n), \end{aligned} \quad (21)$$

for all users $u \in \mathcal{N}$, which can be solved in parallel by the users separately. Therefore, solving \mathcal{P}_2 is synonymous with finding the optimal content delivery node and caching policy, i.e.,

$$\mathcal{P}_3 \arg \max \Omega^u(b_{u,t}^n, x_{u,f}|Z(t)). \quad (22)$$

It is worth noting that users need to make dynamic node selection and caching decisions under constantly changing channel conditions. We transform the joint optimization problem \mathcal{P}_3 of content delivery node selection and cache replacement into a Markov decision process (MDP), as shown below:

State: The state of user at time t can be expressed as $s_u^t = \{Y_u(t), Q_u(t), R_u, h_u(t)\}$. $Y_u(t), Q_u(t)$ denote the energy consumption queue and request queue of user u , respectively. In addition, R_u is the probability of a user establishing a D2D connection with neighbouring user, and $h_u(t)$ is denoted as the channel gain of the user requesting content to other delivery nodes.

Action: After receiving the status s_u^t , select the content sharing transmission node and replace the file. Therefore, the action can be expressed as $a_u^t = \{b_{u,t}^n, I_u^t\}$, where $b_{u,t}^n$ is the indication identifier indicating the node (e.g., including UE's neighbouring nodes and base station). I_u^t is the indicator whether the file in cache list need to be replaced.

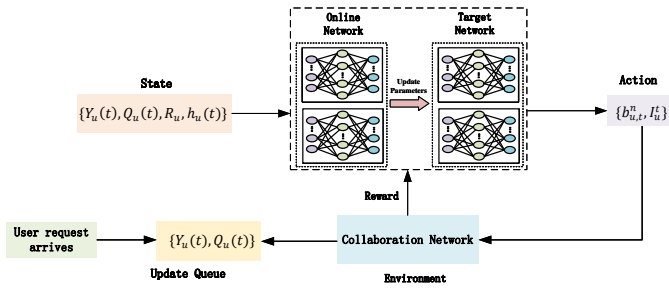


Fig. 3. DDPG architecture for solving the proposed MDP problem.

Reward value of the system: Our goal is to maximize revenue from D2D and content sharing while reducing user request latency and maintaining the stability of the request and energy queues in the system. Therefore, we set the reward as the optimization problem \mathcal{P}_3 , i.e., $r_u = \Omega^u(b_{u,t}^n, x_{u,f} | Z(t))$.

In general, obtaining the optimal action, including content node selection decision $(b_{u,t}^n)^*$ and cache replacement I_u^t , involves exploring a vast number of possible decisions, which can be as high as $2^{(N+M_u)}$ options. This results in significant computational complexity even when N is very small. As a result, a reinforcement learning approach is employed to enable online shared node selection and cache decisions. Specifically, we use the deep deterministic policy gradient (DDPG) reinforcement learning algorithm to dynamically make node selection and cache decisions, thus solving the formulated MDP problem.

Fig. 3 illustrates the architecture for addressing the proposed MDP problem in DDPG. The DDPG algorithm employs two independent deep neural networks (DNNs), which follow the actor-critic paradigm, including an online network and a target network. In each selection period t , the state s_u is sent to the actor network of the online network. To ensure comprehensive exploration of the environment while maintaining a balance between exploration and exploitation, a Gaussian noise vector is incorporated into the action policy function during its output. Following the execution of action a_u , the ensuing reward r_u and the subsequent state $s'_u(t)$ are observed. The observed Transition state $(s_u(t), a_u(t), r_u(t), s'_u(t))$ is saved to the experience pool to facilitate subsequent network learning.

The target network, which acts as a delayed replica of the online network, progressively tracks the acquired knowledge and updates the parameter configuration of the target network model through a process of soft updates. Throughout the training of the system, the agent randomly samples from the experience replay pool. In an effort to ensure the proximity of the Q-value output produced by the critic network to the actual value, the Q-value estimated by the critic network is employed within the target network. Additionally, the mean square error loss function is utilized to guide the training of the actor network. The model is trained Episode times, and in each epoch, the target network implements a soft update mechanism to modify network parameters. The detailed LoRL Algorithm is shown in algorithm 1.

Algorithm 1 Combines Lyapunov optimization and reinforcement learning (LoRL) for solving (\mathcal{P}_3) .

input: Parameters V .

output: Control actions $a_u^t = \{b_{u,t}^n, I_u^t\}$.

initialization: $\{h_u\}_{u=1}^{\mathcal{H}}, Y_u(t)_{u=1}^{\mathcal{H}}, \{Q_u(t)\}_{u=1}^{\mathcal{H}}, Y_u(0) = 0, Q_u(0) = 0$.

for UE- $u \in \mathcal{N}$ **in Parallel:** **do**

Training AutoEncoder model, and then obtain user preferences P_u .

end for

for each epoch **in episode** **do**

for $t = 0, 1, 2, 3, \dots, T$ **do**

for UE- $u \in \mathcal{N}$ **in Parallel:** **do**

User u sends a request file f .

IF $x_{u,f} = 1$:

continue.

Calculate the similarity S_u^v (12) between user u and neighboring users.

Calculate the connection probability between users $R_{u,v}, \frac{R_{u,v}}{\sum_{u \in \mathcal{N}} R_{u,v}} \rightarrow R_u$.

Observer the state $s_u^t = \{Y_u(t), Q_u(t), R_u, h_u(t)\}$.

Observer reward feedback $r_u(t)$ (21), and obtain new observations $s'_u(t)$.

Construct transition $(s_u(t), a_u(t), r_u(t), s'_u(t))$ and then store the transition into experience pool.

Update actor-critic model parameters

end for

end for

end for

C. Complexity Analysis

The computational complexity analysis of the proposed LoRL scheme is as follows: The execution of the LoRL algorithm consists of two parts, namely, joint user association and cache decision, and policy update. Between these two parts, a joint decision action generation is performed in each time frame, while policy updates are less frequent. Therefore, we focus on analyzing the complexity of policy decision generation in each time frame.

Careful observation reveals that within each time slot, the algorithm's complexity includes the computation of similarity for each user (12), the probability of connections between users, and the updating of DDPG model parameters. Specifically, the complexity of computing similarity and user connection probability is $O(2N)$, where N is the number of users.

Additionally, the time complexity of the DDPG algorithm, which mainly consists of initialization, memory replay, and four deep neural networks, is as follows. The state is initialized at the beginning of each training episode, with a time complexity of K . Additionally, both the actor network and the critic network are designed as fully connected networks, assuming the actor network has N_A fully connected layers, and the critic network is composed of N_C fully connected layers. Therefore, the time complexity of DDPG is

$$\begin{aligned}
 & 2 \sum_{i=0}^{N_A-1} u_{a,i} u_{a,i+1} + 2 \sum_{j=0}^{N_C-1} u_{c,j} u_{c,j+1} + v_a u_i + K = \\
 & \mathcal{O} \left(\sum_{i=0}^{N_A-1} u_{a,i} u_{a,i+1} + \sum_{j=0}^{N_C-1} u_{c,j} u_{c,j+1} \right) + \mathcal{O}(K)
 \end{aligned} \tag{23}$$

Based on the description above, the time complexity of the proposed LoRL algorithm is $\mathcal{O}(\sum_{i=0}^{N_A-1} u_{a,i} u_{a,i+1} + \sum_{j=0}^{N_C-1} u_{c,j} u_{c,j+1}) + \mathcal{O}(K) + \mathcal{O}(2N)$.

VI. SIMULATION RESULTS

A. Parameter Setting

In our simulation study, we utilized the MovieLens 1M dataset [36] to model the request behavior in the network. This dataset contains user ratings for a total of 3952 movies. Each record in the dataset includes a user ID, a movie ID, a rating, and a timestamp. Since user ratings are typically provided after viewing, we treated these ratings as request records for our simulation. To calculate user preferences, we divided the dataset into two parts. The period from January 1, 2000 to April 13, 2002 was used as a historical training set to obtain user preferences. The remaining data served as the test set to evaluate the performance of our algorithm. The content database \mathcal{F} consists of the 3952 movies contained in the dataset. In order to reflect the degree of queue backlog in the system network, we set the total number of requests to be 10,000 under different numbers of users. We set the default cache sizes of user nodes and base stations to 40M and 100M, respectively. These cache sizes determine the amount of content that can be stored locally at each node.

TABLE II
MODEL SIMULATION PARAMETERS

| Parameter | Value |
|---------------------|-------|
| learning rate | 0.01 |
| action dim | 400 |
| critic dim | 200 |
| batch size | 128 |
| gamma | 0.99 |
| Exploration rate | 0.1 |
| Hidden layer | 3 |
| Memory Capacity | 2000 |
| Activation function | ReLU |

In our simulation, we modeled the user's movement trajectory using the Random Waypoint model [37], which is a widely applied and proven effective approach in simulating user mobility. The Random Waypoint model has also been commonly used in other studies, particularly in research related to caching strategies [38] [39]. This model generates random coordinate locations for the user at each time slot within a specified region. To determine the range of D2D connections, we assumed that users can establish connections within a physical range limit of 100 units. Users can connect to the base station within a range of 500 units. The path loss in the network was modeled using the formula $36.8 + 36.7 \log(d)$, where d represents the distance between user cache nodes. The small-scale fading was modeled using the unit variance of

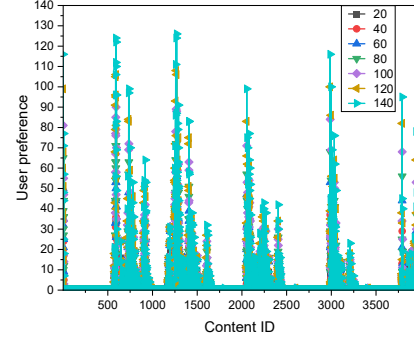


Fig. 4. Variation in request preference at different numbers of users.

Rayleigh fading. Other network parameters included a channel bandwidth of 20 MHz and a background noise level of -95 dBm. In addition, all users run a DRL agent with a three-layer neural network. All these agents use the Adam optimizer with adaptive learning rates to learn their respective training parameters, starting from a learning rate of 10^{-2} . For specific simulation parameter settings of DDPG, please refer to Table II.

B. Content Request Preference Analysis

In our experiment, we simplified the dataset by extracting the user IDs with the highest number of requests. We then selected a specific number of users for simulation purposes. User preference in this context refers to users' focus on requesting certain types of files. To illustrate how user preference changes with different numbers of users, we provide a simple example. In the example, we consider the numbers of users to be [20, 40, 60, 80, 100, 120, 140]. Fig. 4 shows the distribution of requested content IDs at different numbers of users. From the figure, we observe that the content popularity exhibits certain preferences for different numbers of users. The requested content IDs are clustered within specific ranges, i.e., [500-1000], [1250-1500], [2000-2500], [3000-3250], [3725-3952]. This indicates that certain content types or categories are more popular among users, and their preferences can be observed based on the content IDs requested.

C. Baseline Schemes and Performance Metrics

To evaluate the performance of the proposed LoRL algorithm under different parameters, we consider the following baseline scheme:

- 1) **LRU**: Randomly select content sharing nodes within the communicable range, and the earliest stored content will be replaced with new content.
- 2) **MLPLRU** [40]: Randomly select content sharing nodes within the communicable range, and an improved on-line cache replacement strategy based on Pareto Least Recently Used (PLRU) Algorithm and Least Recently Used K (LRU-K) Algorithm.
- 3) **LoRL-with no preference**: LoRL-with no preference refers to the variant of LoRL algorithm that disregards

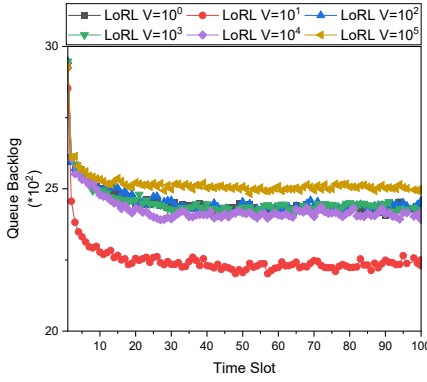


Fig. 5. Dynamics of the time-averaged sum of queue backlogs in LoRL for different weights V during the observation period $t = 1, \dots, 100$ slots.

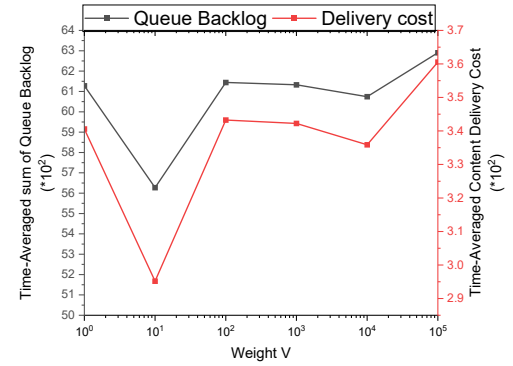


Fig. 7. Time averages of the content transfer cost and the sum of queue backlogs in LoRL as functions of weight V .

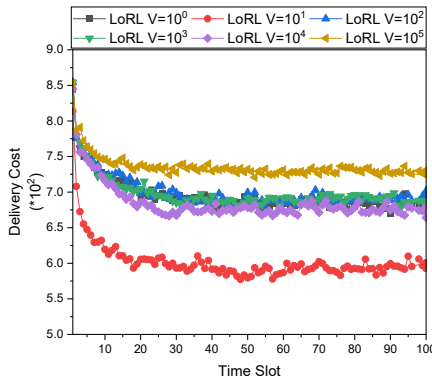


Fig. 6. Dynamics of the time-averaged network cost in LoRL for different weights V during the observation period $t = 1, \dots, 100$ slots.

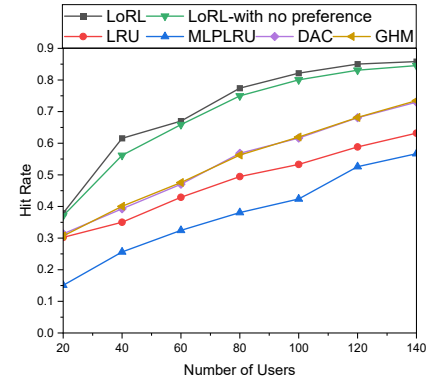


Fig. 8. The cache hit ratio varies with the number of users.

user preferences and solely focuses on the connection probability between users.

- 4) **DAC**: A delay-aware D2D caching (DAC) [41] algorithm with the goal of request latency.
- 5) **GHM**: The greedy heuristic method (GHM) [10] searches within a limited number of user/file pairs to maximize the target D2D gain and delay value.

To evaluate these solutions, we use the following performance metrics: (i) hit rate (satisfied by local cache, D2D sharing, or BS); (ii) average delay; (iii) queue backlog; (iv) network cost; (v) D2D offloading rate.

D. The impact of Weight V

Figures 5-7 verify the analysis results related to Lyapunov optimization established in (17). Figures 5 and 6 show the dynamics of the time-averaged network cost and sum of queue backlogs during different values of weight V , respectively, in LoRL during the period $t = 1, \dots, 100$. It is observed that LoRL converges to stable delivery cost and queue backlog levels around time slot $t = 30$. We also see that when weight $V = 10$, there is lower network cost and queue backlog. Figure 7 illustrates the target value, which is the time average of the sum of transmission cost and queue backlog, as a function of the weight V in Lyapunov optimization. The graph

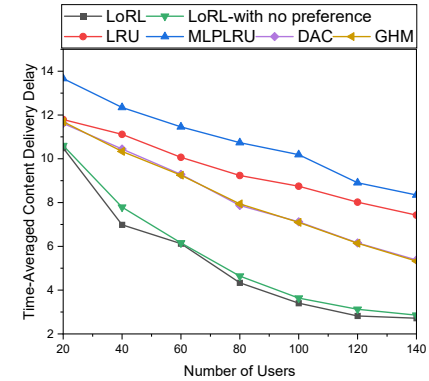


Fig. 9. Time-average transmission delay varies with the number of users.

demonstrates a proportional relationship between the sum of queue backlog and network cost. This can be attributed to the fact that a higher queue backlog prompts the selection of delivery nodes with lower transmission delays to minimize retrieval costs, thereby resulting in an increase in delivery costs.

E. The impact of User Numbers

Results in Figures 8-12 describe the performance of the cache strategy in terms of cache hit ratio, time-averaged

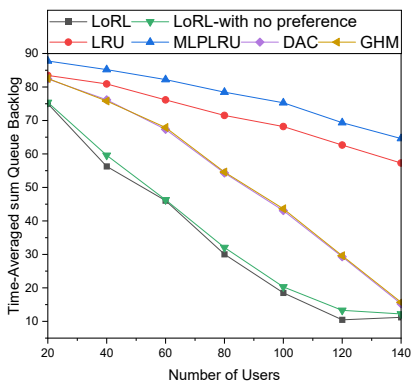


Fig. 10. Time-average sum of queue backlogs with varies with the number of users.

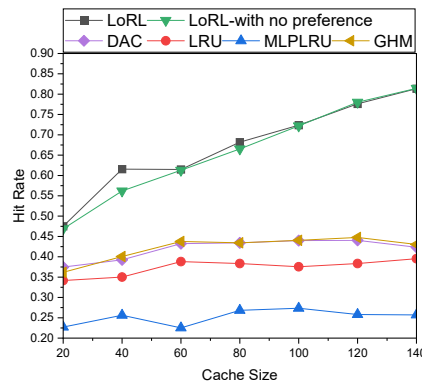


Fig. 13. The cache hit ratio varies with cache size.

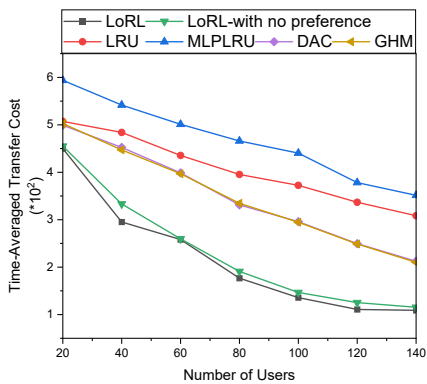


Fig. 11. Time-average network cost varies with the number of users.

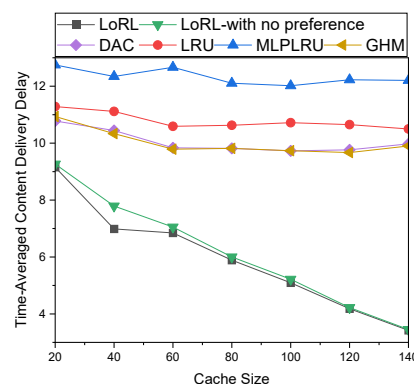


Fig. 14. Time-averaged transmission delay varies with cache size.

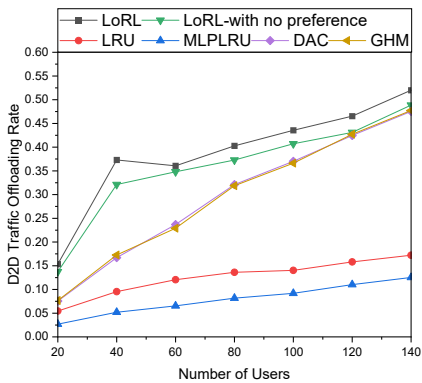


Fig. 12. The Device-to-Device (D2D) offloading rate varies with the number of users.

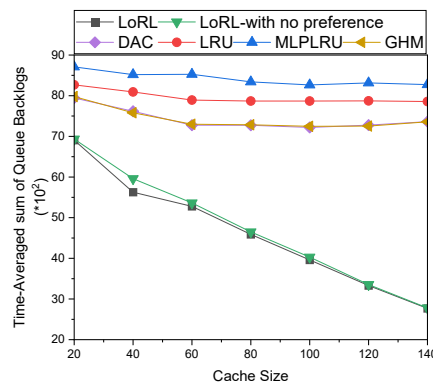


Fig. 15. Time-averaged sum of queue backlogs varies with cache size.

latency, time-averaged delivery cost, and time-averaged queue backlog under different numbers of users. The Figure 8 shows that the higher the number of users, the higher the cache hit rate. As depicted in Figures 9-11, it is evident that with the expansion of the network user scale, the queue backlog, transmission cost, and average request delay decrease accordingly. This can be attributed to the increased opportunities for D2D communication resulting from a larger number of users, thereby enhancing the effectiveness of collaborative caching. It is worth noting that it is observed from Figure 12 that as

the number of users increases, the D2D offloading rate has a slow decreasing trend, which is mainly because our proposed mechanism combines user preference node selection and cache decision, local The cache decision is more biased towards local requests, increasing the probability of local cache hits, thereby reducing the proportion of D2D offloading. Furthermore, it can be seen from these figures that although the proportion of D2D offloading decreases, the proposed mechanism still outperforms other strategies.

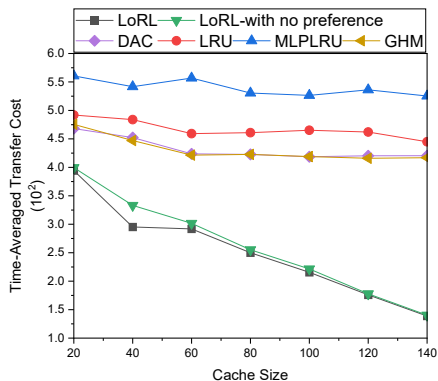


Fig. 16. Time-averaged network cost varies with cache size.

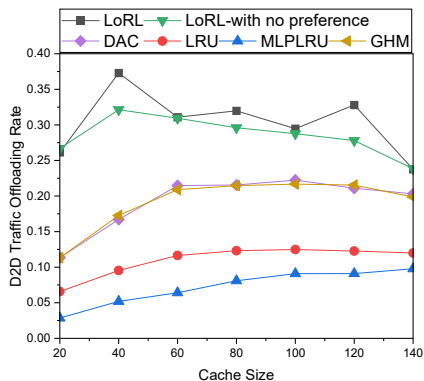


Fig. 17. The Device-to-Device (D2D) offloading rate varies with cache size.

F. The impact of Cache Size

In order to compare the performance of LoRL with other caching strategies, we present the results in Figures 13-17 with respect to cache size. These figures display the hit rate, request latency, time-averaged network cost, and the sum of queue backlog and D2D offload ratio. It is evident that the cache hit ratio improves with larger cache sizes. Conversely, the combined metrics of queue backlogs, transmission delays, and network costs decrease as the user cache size increases. As can be seen from these figures, LRU and MLPLRU are not very effective in minimizing content delivery delay, transmission cost, and queue backlog. The reason is that LRU and MLPLRU do not consider the popularity of file content, and only update the cache according to local requests. Also, these algorithms randomly select shared nodes from within the exchange range. In addition, both the GHM algorithm and the DAC algorithm are greedy node selection algorithms, which can improve the content sharing rate by optimizing the D2D node selection, but it does not consider the local cache decision, so its local cache hit rate is low. We also observed that the DAC and GHM algorithms are inefficient, leading to performance degradation at the expense of queue backlog and network transmission costs. The main reason is that when the number of users is fixed, the cache nodes that users can establish D2D links are limited, so GHM is compared to DAC advantage is not particularly obvious. In particular, LoRL has

certain performance advantages compared to LoRL with no preference, and can better select shared nodes with high D2D gains, making the network have lower queue backlogs and transmission costs.

VII. CONCLUSION

In this paper, we introduce a novel approach for selecting shared nodes and making cache decisions in D2D-assisted heterogeneous collaborative edge computing networks. Our approach involves formulating a joint optimization problem and utilizing a Lyapunov optimization algorithm to decouple the problem. To enable intelligent user association and caching decisions, we propose a content caching algorithm based on deep deterministic policy gradient (DDPG) [35]. The algorithm aims to minimize user request latency while ensuring the stability of request and energy consumption queues in the system. To evaluate the effectiveness of our proposed algorithm, we conduct an extensive study and compare it with five baseline schemes. The results demonstrate that our algorithm surpasses the baseline schemes in terms of average content download latency and system queue stability. In our future work, we intend to delve deeper into the fine-grained characteristics of users and requests. This will enable us to make more informed decisions regarding user association and caching, ultimately enhancing the overall performance of the system.

REFERENCES

- [1] (ITU-R) I R S. Report ITU-R M.2370-0-IMT traffic estimates for the years 2020 to 2030[R]. Geneva: International Telecommunication Union (ITU), 2015.
- [2] Song J, Sheng M, Quek T Q S, et al. Learning-based content caching and sharing for wireless networks[J]. *IEEE Transactions on Communications*, 2017, 65(10): 4309-4324.
- [3] Wang X, Li R, Wang C, et al. Attention-weighted federated deep reinforcement learning for device-to-device assisted heterogeneous collaborative edge caching[J]. *IEEE Journal on Selected Areas in Communications*, 2020, 39(1): 154-169.
- [4] Li L, Zhao G, Blum R S. A survey of caching techniques in cellular networks: Research issues and challenges in content placement and delivery strategies[J]. *IEEE Communications Surveys and Tutorials*, 2018, 20(3): 1710-1732.
- [5] Chen B, Yang C. Caching policy for cache-enabled D2D communications by learning user preference[J]. *IEEE Transactions on Communications*, 2018, 66(12): 6586-6601.
- [6] Zhou H, Wu T, Zhang H, et al. Incentive-driven deep reinforcement learning for content caching and D2D offloading[J]. *IEEE Journal on Selected Areas in Communications*, 2021, 39(8): 2445-2460.
- [7] Ali M S , Coucheney P , Coupechoux M . Learning in Noisy-Potential Games for Resource Allocation in D2D Networks[J]. *IEEE Transactions on Mobile Computing*, 2019, 19(12): 2761-2773.
- [8] Li L, Xu Y, Yin J, et al. Deep reinforcement learning approaches for content caching in cache-enabled D2D networks[J]. *IEEE Internet of Things Journal*, 2019, 7(1): 544-557.
- [9] Lin F P C, Hosseinalipour S, Azam S S, et al. Semi-decentralized federated learning with cooperative D2D local model aggregations[J]. *IEEE Journal on Selected Areas in Communications*, 2021, 39(12): 3851-3869.
- [10] Asheralieva A, Niyato D. Combining contract theory and Lyapunov optimization for content sharing with edge caching and device-to-device communications[J]. *IEEE/ACM Transactions on Networking*, 2020, 28(3): 1213-1226.
- [11] Bi S, Huang L, Wang H, et al. Lyapunov-guided deep reinforcement learning for stable online computation offloading in mobile-edge computing networks[J]. *IEEE Transactions on Wireless Communications*, 2021, 20(11): 7519-7537.

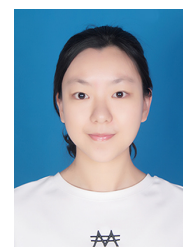
- [12] Yan S, Jiao M, Zhou Y, et al. Machine-learning approach for user association and content placement in fog radio access networks[J]. *IEEE Internet of Things Journal*, 2020, 7(10): 9413-9425.
- [13] Jiang Y, Ma M, Bennis M, et al. User preference learning-based edge caching for fog radio access network[J]. *IEEE Transactions on Communications*, 2018, 67(2): 1268-1283.
- [14] Jiang Y, Wu Y, Zheng F C, et al. Federated Learning Based Content Popularity Prediction in Fog Radio Access Networks[J]. *IEEE Transactions on Wireless Communications*, 2021.
- [15] Chen M, Saad W, Yin C, et al. Echo state networks for proactive caching in cloud-based radio access networks with mobile users[J]. *IEEE Transactions on Wireless Communications*, 2017, 16(6): 3520-3535.
- [16] Hassine N B, Milocco R, Minet P. ARMA based popularity prediction for caching in content delivery networks[C]//2017 *Wireless Days*. IEEE, 2017: 113-120.
- [17] Chen Q, Wang W, Yu F R, et al. Content caching oriented popularity prediction: A weighted clustering approach[J]. *IEEE Transactions on Wireless Communications*, 2020, 20(1): 623-636.
- [18] Wang X, Wang C, Li X, et al. Federated deep reinforcement learning for Internet of Things with decentralized cooperative edge caching[J]. *IEEE Internet of Things Journal*, 2020, 7(10): 9441-9455.
- [19] Li S, Xu J, Van Der Schaar M, et al. Popularity-driven content caching[C]// *IEEE INFOCOM 2016-The 35th Annual IEEE International Conference on Computer Communications*. IEEE, 2016: 1-9.
- [20] J. Kwak, Y. Kim, L. B. Le and S. Chong, "Hybrid Content Caching in 5G Wireless Networks: Cloud Versus Edge Caching," in *IEEE Transactions on Wireless Communications*, vol. 17, no. 5, pp. 3030-3045, May 2018, doi: 10.1109/TWC.2018.2805893
- [21] Qian Y, Wang R, Wu J, et al. Reinforcement learning-based optimal computing and caching in mobile edge network[J]. *IEEE Journal on Selected Areas in Communications*, 2020, 38(10): 2343-2355.
- [22] Somuyiwa S O, György A, Gündüz D. A reinforcement-learning approach to proactive caching in wireless networks[J]. *IEEE Journal on Selected Areas in Communications*, 2018, 36(6): 1331-1344.
- [23] Sadeghi A, Sheikholeslami F, Giannakis G B. Optimal and scalable caching for 5G using reinforcement learning of space-time popularities[J]. *IEEE Journal of Selected Topics in Signal Processing*, 2017, 12(1): 180-190.
- [24] Chen Q, Wang W, Chen W, et al. Cache-enabled multicast content pushing with structured deep learning[J]. *IEEE Journal on Selected Areas in Communications*, 2021, 39(7): 2135-2149.
- [25] Liu S, Zheng C, Huang Y, et al. Distributed Reinforcement Learning for Privacy-Preserving Dynamic Edge Caching[J]. *IEEE Journal on Selected Areas in Communications*, 2022, 40(3): 749-760.
- [26] Huang W, Chen W, Poor H V. Request delay-based pricing for proactive caching: A stackelberg game approach[J]. *IEEE Transactions on Wireless Communications*, 2019, 18(6): 2903-2918.
- [27] Z. Xu et al., "Stable Service Caching in MECs of Hierarchical Service Markets with Uncertain Request Rates," in *IEEE Transactions on Mobile Computing*, doi: 10.1109/TMC.2022.3149870.
- [28] Müller S, Atan O, van der Schaar M, et al. Context-aware proactive content caching with service differentiation in wireless networks[J]. *IEEE Transactions on Wireless Communications*, 2016, 16(2): 1024-1036.
- [29] Saputra Y M, Hoang D T, Nguyen D N, et al. A novel mobile edge network architecture with joint caching-delivering and horizontal cooperation[J]. *IEEE Transactions on Mobile Computing*, 2019, 20(1): 19-31.
- [30] Kwak J, Kim Y, Le L B, et al. Hybrid content caching in 5G wireless networks: Cloud versus edge caching[J]. *IEEE Transactions on Wireless Communications*, 2018, 17(5): 3030-3045.
- [31] Zhang T, Wang Y, Yi W, et al. Two Time-Scale Caching Placement and User Association in Dynamic Cellular Networks[J]. *IEEE Transactions on Communications*, 2022, 70(4): 2561-2574.
- [32] Wang X, Wang C, Li X, et al. Federated deep reinforcement learning for internet of things with decentralized cooperative edge caching[J]. *IEEE Internet of Things Journal*, 2020, 7(10): 9441-9455.
- [33] Yao J, Han T, Ansari N. On mobile edge caching[J]. *IEEE Communications Surveys and Tutorials*, 2019, 21(3): 2525-2553.
- [34] 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 Communications.*, vol. 36, no. 8, pp. 1768–1785, Aug. 2018.
- [35] Lillicrap T P, Hunt J J, Pritzel A, et al. Continuous control with deep reinforcement learning[J]. *arXiv preprint arXiv:1509.02971*, 2015.
- [36] F. M. Harper and J. A. Konstan, "The movielens datasets: History and context," *ACM Trans. Interact. Intell. Syst.*, vol. 5, no. 4, pp. 1–19, Jan. 2015.
- [37] Bettstetter C, Hartenstein H, Perez-Costa X. Stochastic Properties of the Random Waypoint Mobility Model[J]. *Wireless Networks*, 2004, 10(5):p.555-567.
- [38] Banerjee B, Tellambura C. Study of mobility in cache-enabled wireless heterogeneous networks[C]//2017 IEEE wireless communications and networking conference (WCNC). IEEE, 2017: 1-6.
- [39] C. Casetti, C. -F. Chiasserini, M. Fiore, C. -A. La and P. Michiardi, "P2P cache-and-forward mechanisms for mobile ad hoc networks," 2009 IEEE Symposium on Computers and Communications, Sousse, Tunisia, 2009, pp. 386-392, doi: 10.1109/ISCC.2009.5202293.
- [40] F. Miao, D. Chen, and L. Jin, "Multi-level plru cache algorithm for content delivery networks," in 2017 10th *International Symposium on Computational Intelligence and Design (ISCID)*, vol. 1, 2017, pp. 320–323.
- [41] Y. Li, M. C. Gursoy, and S. Velipasalar, "A delay-aware caching algorithm for wireless D2D caching networks," in Proc. IEEE Conf. Comput. Commun. Workshops (INFOCOM WKSHPS), May 2017, pp. 456–461.



Ziyi Teng received the B.S. degree from the North China University of Water Resources and Electric Power, Zhengzhou, China, in 2019. She is currently pursuing the Ph.D. degree with the Beijing University of Technology, Beijing, China. She is a student member of CCF. Her research interests in mobile edge computing.



Juan Fang received the M.S. degree from Jilin University of Technology, Changchun, China in 1997 and her Ph.D. degree from the College of Computer Science, Beijing University of Technology, Beijing, China in 2005. In 1997, she joined the College of Computer Science, Beijing University of Technology. From 2015, she is the professor of Beijing University of Technology. She currently works with the Faculty of Information Technology, Beijing University of Technology, Beijing. Her research interests include high performance computing, heterogeneous intelligent computing and edge computing.



Yaqi Liu received the B.S. degree from the Beijing University of Technology, Beijing, China, in 2020. She is currently pursuing the Ph.D. with the Beijing University of Technology, Beijing, China. She is a student member of CCF. Her research interests in mobile edge computing.