

QECO: A QoE-Oriented Computation Offloading Algorithm based on Deep Reinforcement Learning for Mobile Edge Computing

Iman Rahmati , Hamed Shah-Mansouri , and Ali Movaghar 

Abstract—In the realm of mobile edge computing (MEC), efficient computation task offloading plays a pivotal role in ensuring a seamless quality of experience (QoE) for users. Maintaining a high QoE is paramount in today’s interconnected world, where users demand reliable services. This challenge stands as one of the most primary key factors contributing to handling dynamic and uncertain mobile environment. In this study, we delve into computation offloading in MEC systems, where strict task processing deadlines and energy constraints can adversely affect the system performance. We formulate the computation task offloading problem as a Markov decision process (MDP) to maximize the long-term QoE of each user individually. We propose a distributed QoE-oriented computation offloading (QECO) algorithm based on deep reinforcement learning (DRL) that empowers mobile devices to make their offloading decisions without requiring knowledge of decisions made by other devices. Through numerical studies, we evaluate the performance of QECO. **Simulation results with 50 mobile devices and 5 edge nodes validate that the QECO algorithm is able to efficiently recognize the mobile devices’ demands and exploit the computational resources of edge nodes. Consequently, it can achieve a significant improvement of 34.8% – 75.6% in average QoE compared to several benchmark methods and an existing work.**

Index Terms—Mobile edge computing, computation task offloading, quality of experience, deep reinforcement learning.

I. INTRODUCTION

MOBILE edge computing (MEC) [1] has emerged as a promising technological solution to overcome the challenges faced by mobile devices (MDs) when performing high computational tasks, such as real-time data processing and artificial intelligence applications [2] [3]. In spite of the MDs’ technological advancements, their limited computing power and energy may lead to task drops, processing delays, and an overall poor user experience. By offloading intensive tasks to nearby edge nodes (ENs), MEC effectively empowers computation capability and reduces the delay and energy consumption. This improvement enhances the users’ QoE, especially for time-sensitive computation tasks [4] [5].

Efficient task offloading in MEC is a complex optimization challenge due to the dynamic nature of the network and the variety of MDs and servers involved [6] [7]. In particular, determining the optimal offloading strategy, scheduling the tasks, and selecting the most suitable EN for task offloading

are the main challenges that demand careful consideration. Furthermore, the uncertain requirements and sensitive latency properties of computation tasks pose nontrivial challenges that can significantly impact the computation offloading performance in MEC systems with limited resources.

To cope with the dynamic nature of the network, recent research has proposed several task offloading algorithms using machine learning methods. In particular, reinforcement learning (RL) [8] hold promises to determine optimal decision-making policies by capturing the dynamics of environments and learning strategies for accomplishing long-term objectives. While RL-based methods have demonstrated their effectiveness in handling network dynamics, task offloading still encounters several challenges that require further attention. QoE is a time-varying performance measure that reflects user satisfaction and is not affected only by delay [9], but also by energy consumption [10]. Albeit some existing works have investigated the trade-off between delay and energy consumption [11], they fail to properly address the user demands and fulfill QoE requirements. A more comprehensive approach is required to address the dynamic requirements of individual users in real-time scenarios with multiple MDs and ENs.

In this study, we delve into the computation task offloading problem in MEC systems, where strict task processing deadlines and energy constraints can adversely affect the system’s performance. To mitigate these challenges, we propose a distributed QoE-oriented computation offloading (QECO) algorithm, leveraging deep DRL to enhance decision-making in real-time and uncertain environments. Unlike conventional approaches that require global information, QECO allows MDs to make offloading decisions and select the most suitable EN based on locally available data, including task size, queue status, energy levels, and the historical workload of ENs. Specifically, QECO integrates predictive models to continuously estimate dynamic workloads at ENs, enabling MDs to proactively adjust their offloading strategies in real time, even under high network volatility. With a focus on task completion, delay, and energy consumption, QECO quantifies the QoE for MDs and strikes an optimal balance among these factors. By adopting the appropriate policy based on each MD’s specific requirements at any given time, the QECO algorithm significantly improves the QoE for individual users.

Our main contributions are summarized as follows:

- **Task Offloading Problem in the MEC System: We formulate the task offloading problem in queuing-based MEC as an MDP for time-sensitive tasks, taking into**

I. Rahmati and A. Movaghar are with the Department of Computer Engineering, Sharif University of Technology, Tehran, Iran (email: {iman.rahmati, movaghar}@sharif.edu).

H. Shah-Mansouri is with the Department of Electrical Engineering, Sharif University of Technology, Tehran, Iran (email: hamedsh@sharif.edu).

account the time-varying system environments (e.g., the arrival of new tasks, and the computational requirement of each task). Given task completion, task delay, and energy consumption as the QoE requirements of MDs, this approach concentrates on maximizing the long-term QoE, by effectively utilizing the resources and properly handling the dynamic nature of workload at ENs.

- *DRL-based Offloading Algorithm:* To solve the problem of long-term QoE maximization in highly dynamic mobile environment, we propose QECO algorithm based on DRL that empowers each MD to make offloading decisions independently, without prior knowledge of the other MDs' tasks and offloading models. With a focus on the MD's QoE preference, our approach leverages dueling double deep Q-network (D3QN) [12] to prioritize and strike an appropriate balance between QoE factors. To handle the uncertainty of workloads at ENs, we integrate long short-term memory (LSTM) [13] into D3QN, enabling MDs to forecast ENs' workloads over time and improve the estimation of long-term QoE. Additionally, we provide an analysis of the training convergence and computational complexity of the proposed algorithm.
- *Performance Evaluation:* We conduct comprehensive experiments to evaluate the QECO's performance under various computational workloads, as well as its training convergence with different techniques and hyperparameters. The results demonstrate that our algorithm effectively utilizes the computing resources of MDs and ENs, addressing the dynamic workloads at ENs. It converges more quickly compared to DQN and D3QN methods, while providing a substantial improvement of at least 34.8% in average QoE. This advantage is achieved through prioritizing MD's QoE preferences, resulting in improvements of at least 12.2%, 11.7%, and 10.8% in the average number of completed tasks, task delay, and energy consumption, respectively, when compared to the distributed and collective DRL-based offloading algorithm (DCDRL) [14] and several benchmark methods.

The structure of this paper is as follows. Section II reviews the related work. Section III presents the system model, followed by the problem formulation in Section IV. In Section V, we present the algorithm, while Section VI provides an evaluation of its performance. Finally, we conclude in Section VII.

II. RELATED WORK

In recent years, edge computing has become a popular method for scenarios that need to process massive data with weak devices, such as IoTs [15], Internet of Vehicles (IoVs) [16], [17], and Industrial IoT (IIoT) [18] environment. We explore state-of-the-art RL-based techniques characteristic to discuss how they address task offloading problems in MEC. Table I provides an intuitive comparison of the related works.

A. DRL-based Computation Offloading for MEC

To effectively tackle the challenges of MEC arising from the ever-changing nature of networks, recent research highlights the

effectiveness of RL in adapting to environmental changes and learning optimal strategies. To take advantage of RL methods for computation offloading in MEC, most studies transfer the problem to MDP to model the interactions of the environment. Huang *et al.* in [19], focused on a wireless-powered MEC and proposed a DRL-based approach, capable of attaining near-optimal decisions. This is achieved by selectively considering a compact subset of candidate actions in each iteration. Liu *et al.* in [20] investigated a two-timescale computing offloading and resource allocation problem and proposed a resource coordination algorithm based on multi-agent DRL, which can generate interactive information along with resource decisions. Wu *et al.* in [21] introduced a stochastic game-based resource allocation in the SDN-based MEC network. They used an MDP and proposed a multi-agent RL method to minimize both energy consumption and processing delay. Zhou *et al.* in [22] proposed a Q-learning approach, which is an extension of RL to achieve optimal resource allocation strategies and computation offloading. These works in [19], [20] do not account for delay-sensitive tasks and may overlook critical performance factors in real-time applications. Besides, these studies [19]–[22] primarily investigate simplified MEC networks, focusing on single MEC network models, considering only single-server or single-user scenarios, which may not always reflect real-world scenarios. Moreover, due to the absence of robust underlying resource scheduling systems, issues such as data congestion during transmission and resource allocation can arise, leading to suboptimal system performance. Therefore, more accurate MEC models are needed to address the complexities of dynamic resource scheduling in large-scale scenarios.

Some studies explored task offloading problems in more advanced scenarios. Chen *et al.* in [23] addressed a joint optimization problem involving computation offloading and resource allocation, aiming to reduce both task processing delay and energy consumption across all MDs. To optimize privacy protection and quality of service, authors in [24] investigated the joint computation offloading and power allocation problems for the IIoT network. They modeled the problem as an MDP and proposed a multi-agent deep Q-network (DQN)-based algorithm. Dai *et al.* in [25] introduced the integration of action refinement into DRL and designed an algorithm to concurrently optimize resource allocation and computation offloading. In [26], Zhao *et al.* proposed a computation offloading algorithm based on DRL, which addresses the competition for wireless channels to optimize long-term downlink utility. These works did not take into account the challenges of a dynamic network environment. For instance, in [23], MDs are assumed to have sufficient transmission and computation resources from ENs, which may not always align with real-world conditions. While the works of [24]–[26] investigated the limited resources in MEC, they did not consider the dynamic workload at ENs in decision making processes, imposing significant challenges in overall system performance.

B. Queuing Theory for Efficient Resources Scheduling

To reflecting real-world scenarios more accurately, the authors applied queuing theory to model MEC systems. It helps

TABLE I
COMPARISON OF RELATED WORKS CHARACTERISTIC

Paper	Scenario	Problem	Objective	Model	Method	Drawbacks
[27]	Queueing-based MEC	Online computation offloading	Minimize computation delay and energy consumption	MDP	DDPG + DDQN	Only for a single MEC system model
[22]	heterogeneous MEC	Resource allocation and computation offloading	Minimize energy consumption of the entire system	MDP	DDQN	Only for a single MEC system model
[19]	Wireless-powered MEC	Computation offloading and resource allocations	Maximize computation rate	MIP	DQN	Not consider delay-sensitive tasks
[24] 2	MEC-based IIoT	Joint power allocation and computation offloading	Optimization of privacy protection and quality of service.	MDP	Multi-agent actor-critic network	Not take user's demand into consideration
[26]	MEC	User association and resource allocation	Optimize downlink utility	MDP	Multi-agent D3QN	Not consider system delay and energy consumption
[28]	MEC	Computation offloading and service caching	Minimize average service delay	MDP	hierarchical DRL	Not consider system energy consumption
[21] 1	SDN-based MEC	Stochastic game-based resource allocation	Minimize energy consumption and processing delay	MDP	Multi-agent RL	Not consider tasks with maximum delay tolerance
[25]	MEC	Resource allocation and computation offloading	Minimize energy consumption	MDP	DDPG	Not take user's demand into consideration
[20]	MEC	Resource allocation and computation offloading	Minimize execution cost	POMDP	Multi-agent actor-critic network	Not consider delay-sensitive tasks
[29]	MEC-based IIoT	Resource allocation and computation offloading	Minimize task computation delay and energy consumption	MDP	DRL	Not take user's demand into consideration
[30] +	Queueing-based heterogeneous MEC	Decentralized computation offloading	Optimize system cost and completion rates	MDP	Multi-agent actor-critic network	Not take user's demand into consideration
[31] +	Queueing-based MEC	Computation offloading	Minimize drop rate and computation delay.	MDP	D3QN + LSTM	Not consider system energy consumption
[32] 1	Queueing-based IIoT MEC	Privacy aware computation offloading	Optimization of computation rate and energy consumption	MDP	Actor-critic network	Not take user's demand into consideration
[33] 1	Queueing-based IIoT MEC	Computation offloading	Optimize delay and energy consumption	POMDP	Multi-agent DQN	Not consider tasks with maximum delay tolerance
[34]	Queueing-based MEC	Computation offloading	Optimize drop rate and energy consumption	POMDP	DDPG	Not consider computation delay
[35] -	Queueing-based MEC with EHs	Energy harvesting computation offloading	Minimize drop rate and energy consumption.	MDP	DQN	Not consider computation delay
QECO	Queueing-based MEC	Computation offloading	Maximize the QoE of each MD individually	MDP	D3QN + LSTM	

in modeling the dynamic nature of network, allowing MEC systems to allocate computational and transmission resources more efficiently. To minimize energy consumption, Munir *et al.* [36] developed a semi-distributed approach using a multi-agent RL framework for self-powered MEC. Liao *et al.* in [27] introduced a double DRL algorithm for performing online computation offloading in MEC. This algorithm optimizes transmission power and scheduling of CPU frequency when minimizing both task computation delay and energy consumption. To investigate privacy aware computation offloading problem, Wu *et al.* in [32] proposed a DQN-based method. They transform the problem on MDP to optimize computation rate and energy consumption in a queueing-based IIoT network. In [35], the authors proposed an offloading algorithm using deep Q-learning for wireless-powered IoT devices in MEC systems. This algorithm aims to minimize the task drop rate while the devices solely rely on harvested energy for operation. Huang *et al.* in [34] proposed a DRL-based method based on a partially observable MDP (POMDP), which guarantees the deadlines of real-time tasks while minimizing the total energy consumption of MDs. This algorithm effectively tackles the challenges of dynamic resource allocation in large-scale heterogeneous networks. In above

research [36]–[34], while queueing models were employed in several parts of the system, the offloading process was not approached as a holistic queueing system.

C. Task Offloading Considering EN's Workload

Due to the limited processing capacities of ENs, the processing capacity that an EN allocates to an MD depends on the EN's workload each time. Treating this dynamic workload as part of the environment can negatively impact system performance and requires careful consideration. Specifically, offloaded tasks may experience significant processing delays or be dropped if their deadlines expire. Some existing works have addressed the load levels at the edge nodes and proposed task offloading algorithms. Gong *et al.* in [29] proposed a DRL-based network structure in the industrial IoT (IIoT) systems to jointly optimize task offloading and resource allocation to achieve lower energy consumption and decreased task delay. To optimize delay and energy consumption in an MEC-based IIoT system, Wu *et al.* in [33] proposed a multi-agent DQN-based method based on POMDP to address computation offloading problem. Sun *et al.* in [28] explored both computation offloading and service caching problems in MEC. They formulated an

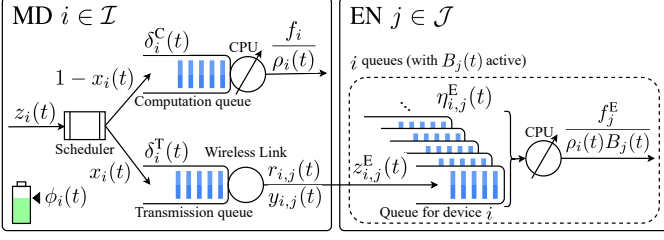


Fig. 1. An illustration of MD $i \in \mathcal{I}$ and EN $j \in \mathcal{J}$ in the MEC system.

optimization problem that aims to minimize the long-term average service delay. They then proposed a hierarchical DRL framework, which effectively handles both problems under heterogeneous resources. Tang *et al.* in [31] investigated the task offloading problem for indivisible and deadline-constrained computational tasks in MEC systems. The authors proposed a distributed DRL-based offloading algorithm designed to handle uncertain workload dynamics at the ENs. In [30], Gao *et al.* introduced an attention-based multi-agent algorithm designed for decentralized computation offloading. However, these algorithms in [29]–[30] require global information of the system and may incur high signaling overhead.

III. SYSTEM MODEL

We investigate a MEC system consisting of a set of MDs denoted by $\mathcal{I} = \{1, 2, \dots, I\}$, along with a set of ENs denoted by $\mathcal{J} = \{1, 2, \dots, J\}$, where I and J represent the number of MDs and ENs, respectively. We regard time as a specific episode containing a series of T time slots denoted by $\mathcal{T} = \{1, 2, \dots, T\}$, each representing a duration of τ seconds. As shown in Fig. 1, we consider two separate queues for each MD to organize tasks for local processing or dispatching to ENs, operating in a first-in-first-out (FIFO) manner. The MD's scheduler is responsible for assigning newly arrived tasks to each of the queues at the beginning of the time slot. On the other hand, we assume that each EN $j \in \mathcal{J}$ consists of I FIFO queues, where each queue corresponds to an MD $i \in \mathcal{I}$. When each task arrives at an EN, it is enqueued in the corresponding MD's queue.

We define $z_i(t)$ as the index assigned to the computation task arriving at MD $i \in \mathcal{I}$ in time slot $t \in \mathcal{T}$. Let $\lambda_i(t)$ denote the size of this task in bits. The size of task $z_i(t)$ is selected randomly from a discrete set $\Lambda = \{\lambda_1, \lambda_2, \dots, \lambda_\theta\}$, where θ represents the number of these values. Hence, $\lambda_i(t) \in \Lambda \cup \{0\}$ to consider the case that no task has arrived. We also denote the task's processing density as $\rho_i(t)$ that indicates the number of CPU cycles required to complete the execution of a unit of the task. Furthermore, we denote the deadline of this task by $\Delta_i(t)$ which is the number of time slots that the task must be completed to avoid being dropped.

We define two binary variables, $x_i(t)$ and $y_{i,j}(t)$ for $i \in \mathcal{I}$ and $j \in \mathcal{J}$ to determine the offloading decision and offloading target, respectively. Specifically, $x_i(t)$ indicates whether task $z_i(t)$ is assigned to the computation queue ($x_i(t) = 0$) or to the transmission queue ($x_i(t) = 1$), and $y_{i,j}(t)$ indicates whether task $z_i(t)$ is offloaded to EN $j \in \mathcal{J}$. If the task is dispatched to EN j , we set $y_{i,j}(t) = 1$; otherwise, $y_{i,j}(t) = 0$.

A. Communication Model

We consider that the tasks in the transmission queue are dispatched to the appropriate ENs via the MD wireless interface. We denote the transmission rate of MD i 's interface when communicating with EN $j \in \mathcal{J}$ in time t as $r_{i,j}(t)$. In time slot $t \in \mathcal{T}$, if task $z_i(t)$ is assigned to the transmission queue for computation offloading, we define $l_i^T(t) \in \mathcal{T}$ to represent the time slot when the task is either dispatched to the EN or dropped. We also define $\delta_i^T(t)$ as the number of time slots that task $z_i(t)$ should wait in the queue before transmission. It should be noted that MD i computes the value of $\delta_i^T(t)$ before making a decision. The value of $\delta_i^T(t)$ is computed as follows:

$$\delta_i^T(t) = \left[\max_{t' \in \{0, 1, \dots, t-1\}} l_i^T(t') - t + 1 \right]^+, \quad (1)$$

where $[\cdot]^+ = \max(0, \cdot)$ and $l_i^T(0) = 0$ for the simplicity of presentation. Note that the value of $\delta_i^T(t)$ only depends on $l_i^T(t)$ for $t' < t$. If MD $i \in \mathcal{I}$ schedules task $z_i(t)$ for dispatching in time slot $t \in \mathcal{T}$, then it will either be dispatched or dropped in time slot $l_i^T(t)$, which is

$$l_i^T(t) = \min \left\{ t + \delta_i^T(t) + \lceil D_i^T(t) \rceil - 1, t + \Delta_i(t) - 1 \right\}, \quad (2)$$

where $D_i^T(t)$ refers to the number of time slots required for the transmission of task $z_i(t)$ from MD $i \in \mathcal{I}$ to EN $j \in \mathcal{J}$. We have

$$D_i^T(t) = \sum_{j \in \mathcal{J}} y_{i,j}(t) \frac{\lambda_i(t)}{r_{i,j}(t)\tau}. \quad (3)$$

Let $E_i^T(t)$ denote the energy consumption of the transmission from MD $i \in \mathcal{I}$ to EN $j \in \mathcal{J}$. We have

$$E_i^T(t) = D_i^T(t) p_i^T(t) \tau, \quad (4)$$

where $p_i^T(t)$ represents the power consumption of the communication link of MD $i \in \mathcal{I}$ in time slot t .

B. Computation Model

The computation tasks can be executed either locally on the MD or on the EN. In this subsection, we provide a detailed explanation of these two cases.

1) *Local Execution:* We model the local execution by a queuing system consisting the computation queue and the MD processor. Let f_i denote the MD i 's processing power (in cycle per second). When task $z_i(t)$ is assigned to the computation queue at the beginning of time slot $t \in \mathcal{T}$, we define $l_i^C(t) \in \mathcal{T}$ as the time slot during which task $z_i(t)$ will either be processed or dropped. If the computation queue is empty, $l_i^C(t) = 0$. Let $\delta_i^C(t)$ denote the number of remaining time slots before processing task $z_i(t)$ in the computation queue. We have:

$$\delta_i^C(t) = \left[\max_{t' \in \{0, 1, \dots, t-1\}} l_i^C(t') - t + 1 \right]^+. \quad (5)$$

In the equation above, the term $\max_{t' \in \{0, 1, \dots, t-1\}} l_i^C(t')$ denotes the time slot at which each existing task in the computation queue, which arrived before time slot t , is either processed or dropped. Consequently, $\delta_i^C(t)$ denotes the number of time slots that task $z_i(t)$ should wait before being processed.

We denote the time slot in which task $z_i(t)$ will be completely processed by $l_i^C(t)$ if it is assigned to the computation queue for local processing in time slot t . We have

$$l_i^C(t) = \min \left\{ t + \delta_i^C(t) + \lceil D_i^C(t) \rceil - 1, t + \Delta_i(t) - 1 \right\}. \quad (6)$$

The task $z_i(t)$ will be immediately dropped if its processing is not completed by the end of the time slot $t + \Delta_i(t) - 1$. In addition, we introduce $D_i^C(t)$ as the number of time slots required to complete the processing of task $z_i(t)$ on MD $i \in \mathcal{I}$. It is given by:

$$D_i^C(t) = \frac{\lambda_i(t)}{f_i \tau / \rho_i(t)}. \quad (7)$$

To compute the MD's energy consumption in the time slot $t \in \mathcal{T}$, we define $E_i^L(t)$ as:

$$E_i^L(t) = D_i^C(t) p_i^C \tau, \quad (8)$$

where $p_i^C = 10^{-27} (f_i)^3$ represents the energy consumption of MD i 's CPU frequency [37].

2) *Edge Execution*: We model the edge execution by the queues associated with MDs deployed at ENs. If computation task $z_i(t')$ is dispatched to EN j in time $t' < t$, we let $z_{i,j}^E(t)$ and $\lambda_{i,j}^E(t)$ (in bits) denote the unique index of the task and the size of the task in the i^{th} queue at EN j . We define $\eta_{i,j}^E(t)$ (in bits) as the length of this queue at the end of time slot $t \in \mathcal{T}$. We refer to a queue as an active queue in a certain time slot if it is not empty. That being said, if at least one task is already in the queue from previous time slots or there is a task arriving at the queue, that queue is active. We define $\mathcal{B}_j(t)$ to denote the set of active queues at EN j in time slot t .

$$\mathcal{B}_j(t) = \left\{ i \mid i \in \mathcal{I}, \lambda_{i,j}^E(t) > 0 \text{ or } \eta_{i,j}^E(t-1) > 0 \right\}. \quad (9)$$

We introduce $b_j(t) \triangleq |\mathcal{B}_j(t)|$ that represents the number of active queues in EN $j \in \mathcal{J}$ in time slot $t \in \mathcal{T}$. In each time slot $t \in \mathcal{T}$, the EN's processing power is divided among its active queues using a generalized processor sharing method [38]. Let variable f_j^E (in cycles per second) represent the computational capacity of EN j . Therefore, EN j can allocate computational capacity of $f_j^E / (\rho_i(t) b_j(t))$ to each MD $i \in \mathcal{B}_j(t)$ during time slot t . To calculate the length of the computation queue for MD $i \in \mathcal{I}$ in EN $j \in \mathcal{J}$, we define $\omega_{i,j}(t)$ (in bits) to represent the number of bits from dropped tasks in that queue at the end of time slot $t \in \mathcal{T}$. The backlog of the queue, referred to as $\eta_{i,j}^E(t)$ is given by:

$$\eta_{i,j}^E(t) = \left[\eta_{i,j}^E(t-1) + \lambda_{i,j}^E(t) - \frac{f_j^E \tau}{\rho_i(t) b_j(t)} - \omega_{i,j}(t) \right]^+. \quad (10)$$

We also define $l_{i,j}^E(t) \in \mathcal{T}$ as the time slot during which the offloaded task $z_{i,j}^E(t)$ is either processed or dropped by EN j . Given the uncertain workload ahead at EN j , neither MD i nor EN j has information about $l_{i,j}^E(t)$ until the corresponding task $z_{i,j}^E(t)$ is either processed or dropped. Let $\hat{l}_{i,j}^E(t)$ represent the time slot at which the execution of task $z_{i,j}^E(t)$ starts. In mathematical terms, for $i \in \mathcal{I}$, $j \in \mathcal{J}$, and $t \in \mathcal{T}$, we have:

$$\hat{l}_{i,j}^E(t) = \max \{ t, \max_{t' \in \{0, 1, \dots, t-1\}} l_{i,j}^E(t') + 1 \}, \quad (11)$$

where $l_{i,j}^E(0) = 0$. Indeed, the initial processing time slot of task $z_{i,j}^E(t)$ at EN should not precede the time slot when the task was enqueued or when the previously arrived tasks were processed or dropped. Therefore, $l_{i,j}^E(t)$ is the time slot that satisfies the following constraints.

$$\sum_{t'=\hat{l}_{i,j}^E(t)}^{l_{i,j}^E(t)} \frac{f_j^E \tau}{\rho_i(t) b_j(t')} \mathbb{1}(i \in \mathcal{B}_j(t')) \geq \lambda_{i,j}^E(t), \quad (12)$$

$$\sum_{t'=\hat{l}_{i,j}^E(t)}^{l_{i,j}^E(t)-1} \frac{f_j^E \tau}{\rho_i(t) b_j(t')} \mathbb{1}(i \in \mathcal{B}_j(t')) < \lambda_{i,j}^E(t), \quad (13)$$

where $\mathbb{1}(z \in \mathbb{Z})$ is the indicator function. In particular, the total processing capacity that EN j allocates to MD i from the time slot $\hat{l}_{i,j}^E(t)$ to the time slot $l_{i,j}^E(t)$ should exceed the size of task $z_{i,j}^E(t)$. Conversely, the total allocated processing capacity from the time slot $\hat{l}_{i,j}^E(t)$ to the time slot $l_{i,j}^E(t) - 1$ should be less than the task's size.

Additionally, we define $D_{i,j}^E(t)$ to represent the quantity of processing time slots allocated to task $z_{i,j}^E(t)$ when executed at EN j . This value is given by:

$$D_{i,j}^E(t) = \frac{\lambda_{i,j}^E(t) \rho_i(t)}{f_j^E \tau / b_j(t)}. \quad (14)$$

We also define $E_{i,j}^E(t)$ as the energy consumption of processing at EN j in time slot t by MD i . This can be calculated as:

$$E_{i,j}^E(t) = \frac{D_{i,j}^E(t) p_j^E \tau}{b_j(t)}, \quad (15)$$

where p_j^E is a constant value which denotes the energy consumption of the EN j 's processor when operating at full capacity.

In addition to the energy consumed by EN j for task processing, we also take into account the energy consumed by the MD i 's user interface in the standby state while waiting for task completion at the EN j . We define $E_{i,j}^I(t)$ as the energy consumption associated with the user interface of MD $i \in \mathcal{I}$, which is given by

$$E_i^I(t) = \sum_{j \in \mathcal{J}} y_{i,j}(t) D_{i,j}^E(t) p_i^I \tau, \quad (16)$$

where p_i^I is the standby energy consumption of MD $i \in \mathcal{I}$.

$$E_i^O(t) = E_i^T(t) + \sum_{j \in \mathcal{J}} y_{i,j}(t) E_{i,j}^E(t) + E_i^I(t). \quad (17)$$

IV. TASK OFFLOADING PROBLEM FORMULATION

Based on the introduced system model, we present the computation task offloading problem in this section. Our primary goal is to enhance each MD's QoE individually by taking the dynamic demands of MDs into account. To achieve this, we approach the optimization problem as an MDP, aiming to maximize the MD's QoE by striking a balance among key QoE factors, including task completion, task delay, and energy consumption. To prioritize QoE factors, we utilize the MD's **energy level**, which plays a crucial role in decision-making. Specifically, when an MD observes its state (e.g. task size,

queue details, and [energy status](#)) and encounters a newly arrived task, it selects an appropriate action for that task. The selected action, based on the observed state, will result in enhanced QoE. Each MD strives to maximize its long-term QoE by optimizing the policy mapping from states to actions. In what follows, we first present the state space, action space, and QoE function, respectively. We then formulate the QoE maximization problem for each MD.

A. State Space

A state in our MDP represents a conceptual space that comprehensively describes the state of an MD facing the environment. We represent the MD i 's state in time slot t as vector $\mathbf{s}_i(t)$ that includes the newly arrived task size, the queues information, the MD's [energy level](#), and the workload history at the ENs. The MD observes this vector at the beginning of each time slot. The vector $\mathbf{s}_i(t)$ is defined as follows:

$$\mathbf{s}_i(t) = (\lambda_i(t), \delta_i^C(t), \delta_i^T(t), \boldsymbol{\eta}_i^E(t-1), \phi_i(t), \mathcal{H}(t)), \quad (18)$$

where vector $\boldsymbol{\eta}_i^E(t-1) = (\eta_{i,j}^E(t-1))_{j \in \mathcal{J}}$ represents the queues length of MD i in ENs at the previous time slot and is computed by the MD according to (10). Let $\phi_i(t)$ denote the [energy level](#) of MD i in time slot t . Considering the power modes of [some](#) real mobile device, $\phi_i(t)$ is derived from the discrete set $\Phi = \{\phi_1, \phi_2, \phi_3\}$, corresponding to ultra power-saving, power-saving, and performance modes, respectively.

In addition, to predict future EN workloads, we define the matrix $\mathcal{H}(t)$ as historical data, indicating the number of active queues for all ENs. This data is recorded over T^s time slots, ranging from $t-T^s$ to $t-1$, in $T^s \times J$ matrix. For EN j workload history at i^{th} time slot from $T^s - t$, we define $h_{i,j}(t)$ as:

$$h_{i,j}(t) = b_j(t - T^s + i - 1). \quad (19)$$

EN $j \in \mathcal{J}$ broadcasts $b_j(t)$ at the end of each time slot.

We define vector \mathcal{S} as the discrete and finite state space for each MD. The size of the set \mathcal{S} is given by $\Lambda \times T^2 \times \mathcal{U} \times 3 \times I^{T^s \times J}$, where \mathcal{U} is the set of available queue length values at an EN over T time slots.

B. Action Space

The action space represents the agent's behavior and the decisions. In this context, we define $\mathbf{a}_i(t)$ to denote the action taken by MD $i \in \mathcal{I}$ in time slot $t \in \mathcal{T}$. These actions involve two decisions, (a) Offloading decision to determine whether or not to offload the task, and (b) Offloading target to determine the EN to send the offloaded tasks. Thus, the action of MD i in time slot t can be concisely expressed as the following action tuple:

$$\mathbf{a}_i(t) = (x_i(t), \mathbf{y}_i(t)), \quad (20)$$

where vector $\mathbf{y}_i(t) = (y_{i,j}(t))_{j \in \mathcal{J}}$ represents the selected EN for offloading this task. In Section [V-B](#), we will discuss about the size of this action space.

C. QoE Function

The QoE function evaluates the influence of agent's actions based on its personalized QoE requirement, by taking several key performance factors into account. It incorporates a multi-dimensional and adaptive structure that establishes an appropriate balance among these QoE factors according to the MD's preference, such as individual preferences for reducing task delays or energy consumption, which may vary over time. To reflect the agent's personalized QoE requirements, an adaptive trade-off between task completion and cost reduction is designed, where a weighted sum between completion delay and energy consumption plays as an associated cost. The MDs adjusted the weights dynamically to define the significance of each factor based on its energy mode—performance, power-saving, or ultra power-saving. Based on the system model, we first formulate the calculation of task delay and energy consumption and then represent the associated cost and QoE function accordingly.

Given the selected action $\mathbf{a}_i(t)$ in the observed state $\mathbf{s}_i(t)$, we represent $\mathcal{D}_i(\mathbf{s}_i(t), \mathbf{a}_i(t))$ as the delay of task $z_i(t)$, which indicates the number of time slots from time slot t to the time slot in which task $z_i(t)$ is processed. It is calculated by:

$$\mathcal{D}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)) = (1 - x_i(t)) \left(l_i^C(t) - t + 1 \right) + x_i(t) \left(\sum_{\mathcal{J}} \sum_{t'=t}^T \mathbb{1}(z_{i,j}^E(t') = z_i(t)) y_{i,j}(t) l_{i,j}^E(t') - t + 1 \right), \quad (21)$$

where $\mathcal{D}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)) = 0$ if task $z_i(t)$ is dropped. Accordingly, the energy consumption for task $z_i(t)$ under action $\mathbf{a}_i(t)$ in observed state $\mathbf{s}_i(t)$ is denoted as $\mathcal{E}_i(\mathbf{s}_i(t), \mathbf{a}_i(t))$, given by:

$$\mathcal{E}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)) = (1 - x_i(t)) E_i^L(t) + x_i(t) \left(\sum_{\mathcal{J}} \sum_{t'=t}^T \mathbb{1}(z_{i,j}^E(t') = z_i(t)) y_{i,j}(t) E_i^O(t) \right). \quad (22)$$

To define associated cost, we use a weighted sum between task delay $\mathcal{D}_i(\mathbf{s}_i(t), \mathbf{a}_i(t))$ and energy consumption $\mathcal{E}_i(\mathbf{s}_i(t), \mathbf{a}_i(t))$, where the weights adjust dynamically depending on the MD's energy levels to reflect the MD's preference for each factor. Given the delay and energy consumption of task $z_i(t)$, we define $\mathcal{C}_i(\mathbf{s}_i(t), \mathbf{a}_i(t))$ that denotes the associated cost of task $z_i(t)$ given the action $\mathbf{a}_i(t)$ in the state $\mathbf{s}_i(t)$.

$$\mathcal{C}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)) = \phi_i(t) \mathcal{D}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)) + (1 - \phi_i(t)) \mathcal{E}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)), \quad (23)$$

where $\phi_i(t)$ represents the MD i 's [energy level](#). When the MD is operating in performance mode, the MD's preference is to minimize task delays, thus the delay contributes more to the cost than energy consumption. In contrast, when it switches to ultra power-saving mode, the main attention is directed toward reducing energy consumption and the energy consumption contributes more to the cost than delay.

Finally, we define an adaptive trade-off between task completion and cost reduction as the QoE associated with task $z_i(t)$. The QoE is affected by the increase in the number of completed tasks given by the task's complexity and reducing

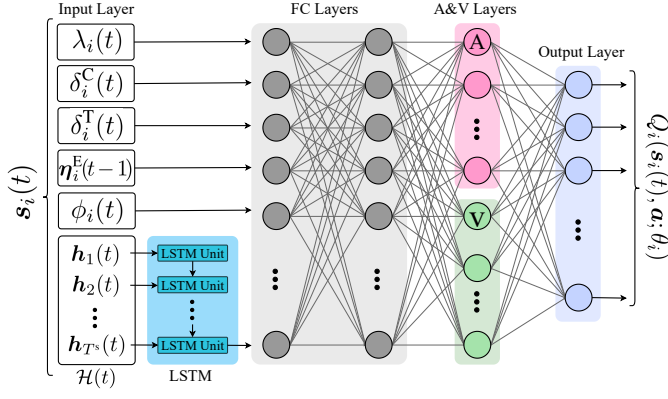


Fig. 2. The neural network of MD $i \in \mathcal{I}$, which characterize the Q-value of each action $a \in \mathcal{A}$ under state $s_i(t) \in \mathcal{S}$.

weighted cost. Given the selected action $a_i(t)$ and the observed state $s_i(t)$, we define The QoE function as follows,

$$q_i(s_i(t), a_i(t)) = \mathcal{R} \lambda_i(t) \rho_i(t) - C_i(s_i(t), a_i(t)) \quad (24)$$

where $\mathcal{R} > 0$ represents a constant reward for task completion. If $z_i(t) = 0$, then $q_i(s_i(t), a_i(t)) = 0$. Throughout the rest of this paper, we adopt the shortened notation $q_i(t)$ to represent $q_i(s_i(t), a_i(t))$.

D. Problem Formulation

We define the task offloading policy for MD $i \in \mathcal{I}$ as a mapping from its state to its corresponding action, denoted by i.e., $\pi_i : \mathcal{S} \rightarrow \mathcal{A}$. Especially, MD i determines an action $a_i(t) \in \mathcal{A}$, according to policy π_i given the observed environment state $s_i(t) \in \mathcal{S}$. The MD aims to find its optimal policy π_i^* which maximizes the long-term QoE,

$$\pi_i^* = \arg \max_{\pi_i} \mathbb{E} \left[\sum_{t \in \mathcal{T}} \gamma^{t-1} q_i(t) \middle| \pi_i \right], \quad (25)$$

where $\gamma \in (0, 1]$ is a discount factor and determines the balance between instant QoE and long-term QoE. As γ approaches 0, the MD prioritizes QoE within the current time slot exclusively. Conversely, as γ approaches 1, the MD increasingly factors in the cumulative long-term QoE. The expectation $\mathbb{E}[\cdot]$ is taken into consideration of the time-varying system environments. Solving the optimization problem in (25) is particularly challenging due to the dynamic nature of the network. To address this challenge, we introduce a DRL-based offloading algorithm to learn the mapping between each state-action pair and their long-term QoE.

V. DRL-BASED OFFLOADING ALGORITHM

We now present QECO algorithm so as to address the distributed offloading decision-making of MDs. The aim is to empower MDs to identify the most efficient action that maximizes their long-term QoE. In the following, we introduce a neural network that characterizes the MD's state-action Q-values mapping, followed by a description of the information exchange between the MDs and ENs.

A. DQN-based Approach

We utilize the DQN technique to find the mapping between each state-action pair to Q-values in the formulated MDP. As shown in Fig. 2, each MD $i \in \mathcal{I}$ is equipped with a neural network comprising six layers. These layers include an input layer, an LSTM layer, two dense layers, an advantage-value (A&V) layer, and an output layer. The parameter vector θ_i of MD i 's neural network is defined to maintain the connection weights and neuron biases across all layers. For MD $i \in \mathcal{I}$, we utilize the state information as the input of neural network. The state information $\lambda_i(t)$, $\delta_i^C(t)$, $\delta_i^T(t)$, $\phi_i(t)$, and $\eta_i^E(t-1)$ are directly passed to the dense layer, while the state information $\mathcal{H}(t)$ is first supplied to the LSTM layer and then the resulting output is sent to the dense layer. The role and responsibilities of each layer are detailed as follows.

1) *Predicting Workloads at ENs*: In order to capture the dynamic behavior of workloads at the ENs, we employ the LSTM network [13]. This network maintains a memory state $\mathcal{H}(t)$ that evolves over time, enabling the neural network to predict future workloads at the ENs based on historical data. By taking the matrix $\mathcal{H}(t)$ as an input, the LSTM network learns the patterns of workload dynamics. The architecture of the LSTM consists of T^s units, each equipped with a set of hidden neurons, and it processes individual rows of the matrix $\mathcal{H}(t)$ sequentially. Through this interconnected design, MD tracks the variations in sequences from $h_1(t)$ to $h_{T^s}(t)$, where vector $h_i(t) = (h_{i,j}(t))_{j \in \mathcal{J}}$, thereby revealing workload fluctuations at the ENs across different time slots. The final LSTM unit produces an output that encapsulates the anticipated workload dynamics, and is then connected to the subsequent layer neurons for further learning.

2) *State-Action Q-Value Mapping*: The pair of dual dense layers plays a crucial role in learning the mapping of Q-values from the current state and the learned load dynamics to the corresponding actions. The dense layers consist of a cluster of neurons that employ rectified linear units (ReLUs) as their activation functions. In the initial dense layer, connections are established from the neurons in the input layer and the LSTM layer to each neuron in the dense layer. The resulting output of a neuron in the dense layer is connected to each neuron in the subsequent dense layer. In the second layer, the outputs from each neuron establish connections with all neurons in the A&V layers.

3) *Dueling-DQN Approach for Q-Value Estimation*: In the neural network architecture, the A&V layer and the output layer incorporate the principles of the dueling-DQN [12] to compute action Q-values. The fundamental concept of dueling-DQN involves two separate learning components: one for action-advantage values and another for state-value. This approach enhances Q-value estimation by separately evaluating the long-term QoE attributed to states and actions.

The A&V layer consists of two distinct dense networks referred to as network A and network V. Network A's role is to learn the action-advantage value for each action, while network V focuses on learning the state-value. For an MD $i \in \mathcal{I}$, we define $V_i(s_i(t); \theta_i)$ and $A_i(s_i(t), a; \theta_i)$ to denote the state-value and the action-advantage value of action $a \in \mathcal{A}$ under state $s_i(t) \in \mathcal{S}$, respectively. The parameter θ_i is responsible

Algorithm 1 QECO Algorithm (Offloading Decision)

Input: state space \mathcal{S} , action space \mathcal{A}
Output: MD $i \in \mathcal{I}$ experience $(\mathbf{s}_i(t), \mathbf{a}_i(t), \mathbf{q}_i(t), \mathbf{s}_i(t+1))$

```

1: for episode 1 to  $N^{\text{ep}}$  do
2:   Initialize  $\mathbf{s}_i(1)$ 
3:   for time slot  $t \in \mathcal{T}$  do
4:     if MD  $i$  receives a new task  $z_i(t)$  then
5:       Send an UpdateRequest to EN  $j_i$ ;
6:       Receive network parameter vector  $\theta_i^E$ ;
7:       Select action  $\mathbf{a}_i(t)$  based on (27);
8:     end if
9:     Observe a set of QoEs  $\{\mathbf{q}_i(t'), t' \in \mathcal{F}_i^t\}$ ;
10:    Observe the next state  $\mathbf{s}_i(t+1)$ ;
11:    for each task  $z_i(t')$  where  $t' \in \mathcal{F}_i^t$  do
12:      Send  $(\mathbf{s}_i(t'), \mathbf{a}_i(t'), \mathbf{q}_i(t'), \mathbf{s}_i(t'+1))$  to EN  $j_i$ ;
13:    end for
14:  end for
15: end for

```

for determining these values, and it can be adjusted when training the QECO algorithm.

For an MD $i \in \mathcal{I}$, the A&V layer and the output layer collectively determine $Q_i(\mathbf{s}_i(t), \mathbf{a}; \theta_i)$, representing the resulting Q-value under action $\mathbf{a} \in \mathcal{A}$ and state $\mathbf{s}_i(t) \in \mathcal{S}$, as follows:

$$Q_i(\mathbf{s}_i(t), \mathbf{a}; \theta_i) = V_i(\mathbf{s}_i(t); \theta_i) + \left(A_i(\mathbf{s}_i(t), \mathbf{a}; \theta_i) - \frac{1}{|\mathcal{A}|} \sum_{\mathbf{a}' \in \mathcal{A}} (A_i(\mathbf{s}_i(t), \mathbf{a}'; \theta_i)) \right), \quad (26)$$

where θ_i establishes a functional relationship that maps Q-values to pairs of state-action.

B. QoE-Oriented DRL-Based Algorithm

The QECO algorithm is meticulously designed to optimize the allocation of computational tasks between MDs and ENs. Since the training of neural networks imposes an extensive computational workload on MDs, we enable MDs to utilize ENs for training their neural networks, effectively reducing their computational workload. For each MD $i \in \mathcal{I}$, there is an associated EN, denoted as EN $j_i \in \mathcal{J}$, which assists in the training process. This EN possesses the highest transmission capacity among all ENs. We define $\mathcal{I}_j \subset \mathcal{I}$ as the set of MDs for which training is executed by EN $j \in \mathcal{J}$, i.e. $\mathcal{I}_j = \{i \in \mathcal{I} | j_i = j\}$. This approach is feasible due to the minimal information exchange and processing requirements for training compared to MD's tasks. The algorithms to be executed at MD $i \in \mathcal{I}$ and EN $j \in \mathcal{J}$ are given in Algorithms 1 and 2, respectively. The core concept involves training neural networks with MD experiences (i.e., state, action, QoE, next state) to map Q-values to each state-action pair. This mapping allows MD to identify the action in the observed state with the highest Q-value and maximize its long-term QoE.

In detail, EN $j \in \mathcal{J}$ maintains a replay buffer denotes as \mathcal{M}_j with two neural networks for MD i : Net_i^E , denoting the evaluation network, and Net_i^T , denoting the target network, which have the same neural network architecture. However,

they possess distinct parameter vectors θ_i^E and θ_i^T , respectively. Their Q-values are represented by $Q_i^E(\mathbf{s}_i(t), \mathbf{a}; \theta_i^E)$ and $Q_i^T(\mathbf{s}_i(t), \mathbf{a}; \theta_i^T)$ for MD $i \in \mathcal{I}_j$, respectively, associating the action $\mathbf{a} \in \mathcal{A}$ under the state $\mathbf{s}_i(t) \in \mathcal{S}$. The replay buffer records the observed experience $(\mathbf{s}_i(t), \mathbf{a}_i(t), \mathbf{q}_i(t), \mathbf{s}_i(t+1))$ of MD i . Moreover, Net_i^E is responsible for action selection, while Net_i^T characterizes the target Q-values, which represent the estimated long-term QoE resulting from an action in the observed state. The target Q-value serves as the reference for updating the network parameter vector θ_i^E . This update occurs through the minimization of disparities between the Q-values under Net_i^E and Net_i^T . In the following, we introduce the offloading decision algorithm of MD $i \in \mathcal{I}$ and the training process algorithm running in EN $j \in \mathcal{J}$.

1) *Offloading Decision Algorithm at MD $i \in \mathcal{I}$* : We analyze a series of episodes, where N^{ep} denotes the number of them. At the beginning of each episode, if MD $i \in \mathcal{I}$ receives a new task $z_i(t)$, it initializes the state $\mathbf{s}_i(1)$ and sends an *UpdateRequest* to EN j_i . After receiving the requested vector θ_i^E of Net_i^E from EN j_i , MD i chooses the following action for task $z_i(t)$.

$$\mathbf{a}_i(t) = \begin{cases} \arg \max_{\mathbf{a} \in \mathcal{A}} Q_i^E(\mathbf{s}_i(t), \mathbf{a}; \theta_i^E), & \text{w.p. } 1 - \epsilon, \\ \text{pick a random action from } \mathcal{A}, & \text{w.p. } \epsilon, \end{cases} \quad (27)$$

where w.p. stands for with probability, and ϵ represents the random exploration probability. The value of $Q_i^E(\mathbf{s}_i(t), \mathbf{a}; \theta_i^E)$ indicates the Q-value under the parameter θ_i^E of the neural network Net_i^E . Specifically, the MD with a probability of $1 - \epsilon$ selects the action associated with the highest Q-value under Net_i^E in the observed state $\mathbf{s}_i(t)$.

In the next time slot $t+1$, MD i observes the state $\mathbf{s}_i(t+1)$. However, due to the potential for tasks to extend across multiple time slots, QoE $\mathbf{q}_i(t)$ associated with task $z_i(t)$ may not be observable in time slot $t+1$. On the other hand, MD i may observe a group of QoEs associated with some tasks $z_i(t')$ in time slots $t' \leq t$. For each MD i , we define the set $\mathcal{F}_i^t \subset \mathcal{T}$ to denote the time slots during which each arriving task $z_i(t')$ is either processed or dropped in time slot t , as given by:

$$\mathcal{F}_i^t = \left\{ t' \left| t' \leq t, \lambda_i(t') > 0, (1 - x_i(t')) l_{i,j}^C(t') + x_i(t') \sum_{j \in \mathcal{J}} \sum_{n=t'}^t \mathbb{1}(z_{i,j}^E(n) = z_i(t')) y_{i,j}(x) l_{i,j}^E(n) = t \right. \right\}.$$

Therefore, MD i observes a set of QoEs $\{\mathbf{q}_i(t') | t' \in \mathcal{F}_i^t\}$ at the beginning of time slot $t+1$, where the set \mathcal{F}_i^t for some $i \in \mathcal{I}$ can be empty. Subsequently, MD i sends its experience $(\mathbf{s}_i(t), \mathbf{a}_i(t), \mathbf{q}_i(t), \mathbf{s}_i(t+1))$ to EN j_i for each task $z_i(t')$ in $t' \in \mathcal{F}_i^t$.

2) *Training Process Algorithm at EN $j \in \mathcal{J}$* : Upon initializing the replay buffer \mathcal{M}_j with the neural networks Net_i^E and Net_i^T for each MD $i \in \mathcal{I}_j$, EN $j \in \mathcal{J}$ waits for messages from the MDs in the set \mathcal{I}_j . When EN j receives an *UpdateRequest* signal from an MD $i \in \mathcal{I}_j$, it responds by transmitting the updated parameter vector θ_i^E , obtained from Net_i^E , back to MD i . On the other side, if EN j receives an experience $(\mathbf{s}_i(t), \mathbf{a}_i(t), \mathbf{q}_i(t), \mathbf{s}_i(t+1))$ from MD $i \in \mathcal{I}_j$, the EN stores this experience in the replay buffer \mathcal{M}_j associated with that MD.

Algorithm 2 QECO Algorithm (Training Process)

Input: experience $(s_i(t), a_i(t), q_i(t), s_i(t+1))$ from MD $i \in \mathcal{I}$
Output: The parameters vector θ_i^E

```

1: Initialize replay buffer  $\mathcal{M}_i$  for each MD  $i \in \mathcal{I}_j$ ;
2: Initialize  $Net_i^E$  and  $Net_i^T$  with random parameters  $\theta_i^E$  and  $\theta_i^T$  respectively, for each MD  $i \in \mathcal{I}_j$ ;
3: Set Count := 0
4: while True do  $\triangleright$  infinite loop
5:   if receive an UpdateRequest from MD  $i \in \mathcal{I}_j$  then
6:     Send  $\theta_i^E$  to MD  $i \in \mathcal{I}_j$ ;
7:   end if
8:   if an experience  $(s_i(t), a_i(t), q_i(t), s_i(t+1))$  is received
9:   from MD  $i \in \mathcal{I}_j$  then
10:    Store  $(s_i(t'), a_i(t'), q_i(t'), s_i(t'+1))$  in  $\mathcal{M}_i$ ;
11:    Get a collection of experiences  $\mathcal{I}$  from  $\mathcal{M}_i$ ;
12:    for each experience  $i \in \mathcal{I}$  do
13:      Get experience  $(s_i(n), a_i(n), q_i(n), s_i(n+1))$ ;
14:      Generate  $\hat{Q}_{i,n}^T$  according to (28);
15:    end for
16:    Set vector  $\hat{\mathbf{Q}}_i^T := (\hat{Q}_{i,n}^T)_{n \in \mathcal{N}}$ ;
17:    Update  $\theta_i^E$  to minimize  $L(\theta_i^E, \hat{\mathbf{Q}}_i^T)$  in (30);
18:    Count := Count + 1;
19:    if mod(Count, ReplaceThreshold) = 0 then
20:       $\theta_i^T := \theta_i^E$ ;
21:    end if
22:  end if
23: end while

```

The EN randomly selects a sample collection of experiences from the replay buffer, denoted as \mathcal{N} . For each experience $n \in \mathcal{N}$, it calculates the value of $\hat{Q}_{i,n}^T$. This value represents the QoE in experience n and includes a discounted Q-value of the action anticipated to be taken in the subsequent state of experience n , according to the network Net_i^T , given by

$$\hat{Q}_{i,n}^T = q_i(n) + \gamma Q_i^T(s_i(n+1), \tilde{a}_n; \theta_i^T), \quad (28)$$

where \tilde{a}_n denotes the optimal action for the state $s_i(n+1)$ based on its highest Q-value under Net_i^E , as given by:

$$\tilde{a}_n = \arg \max_{a \in \mathcal{A}} Q_i^E(s_i(n+1), a; \theta_i^E). \quad (29)$$

In particular, regarding experience n , the target-Q value $\hat{Q}_{i,n}^T$ represents the long-term QoE for action $a_i(n)$ under state $s_i(n)$. This value corresponds to the QoE observed in experience n , as well as the approximate expected upcoming QoE. Based on the previous sample experiences \mathcal{N} , the EN computes the vector $\hat{\mathbf{Q}}_i^T = (\hat{Q}_{i,n}^T)_{n \in \mathcal{N}}$ and trains the MD's neural network (in steps 11-21 of Algorithm 2) to keep parameter vector θ_i^E in Net_i^E update for the MD's *UpdateRequest*. The key idea of updating Net_i^E is to minimize the disparity in Q-values between Net_i^E and Net_i^T , as indicated by the following loss function:

$$L(\theta_i^E, \hat{\mathbf{Q}}_i^T) = \frac{1}{|\mathcal{N}|} \sum_{n \in \mathcal{N}} \left(Q_i^E(s_i(n), a_i(n); \theta_i^E) - \hat{Q}_{i,n}^T \right)^2. \quad (30)$$

In every *ReplaceThreshold* iterations, the update of Net_i^T will involve duplicating the parameters from Net_i^E ($\theta_i^T = \theta_i^E$). The

objective is to consistently update the network parameter θ_i^T in Net_i^T , which enhances the approximation of the long-term QoE when computing the target Q-values in (28).

3) *Computational Complexity*: The computational complexity of the QECO algorithm is determined by the number of experiences required to discover the optimal offloading policy. Each experience involves backpropagation for training, which has a computational complexity of $\mathcal{O}(C)$, where C represents the number of multiplication operations in the neural network. During each training round triggered by the arrival of a new task, a sample collection of experiences of size $|\mathcal{N}|$ is utilized from the replay buffer. Since the training process encompasses N^{ep} episodes and there are K expected tasks in each episode, the computational complexity of the proposed algorithm is $\mathcal{O}(N^{\text{ep}} K |\mathcal{N}| C)$, which is polynomial. Given the integration of neural networks for function approximation, the convergence guarantee of the DRL algorithm remains an open problem. In this work, we will empirically evaluate the convergence of the proposed algorithm in Section VI-B.

VI. PERFORMANCE EVALUATION

In this section, we first present the simulation setup and training configuration. We then illustrate the convergence of the proposed DRL-based QECO algorithm under different optimization methods and hyperparameters, and evaluate its performance in comparison to three baseline schemes in addition to the existing work [14].

A. Simulation Setup

We consider a MEC environment with 50 MDs and 5 ENs, similar to [31]. We also follow the model presented in [22] to determine the energy consumption. All the parameters are given in Table II. To train the MDs' neural networks, we adopt a scenario comprising 1000 episodes. Each episode contains 100 time slots, each of length 0.1 second. The QECO algorithm incorporates real-time experience into its training process to continuously enhance the offloading strategy. Specifically, we employ a batch size of 16, maintain a fixed learning rate of 0.001, and set the discount factor γ to 0.9. The probability of random exploration gradually decreases from an initial value 1, progressively approaching 0.01, all of which is facilitated by an RMSProp optimizer. Detailed settings regarding the neural networks and simulation can be found at our github repository.¹

We use the following methods as benchmarks.

- 1) **Local Computing (LC)**: The MDs execute all of their computation tasks using their own computing capacity.
- 2) **Full Offloading (FO)**: Each MD dispatches all of its computation tasks while choosing the offloading target randomly.
- 3) **Random Decision (RD)**: In this approach, when an MD receives a new task, it randomly makes the offloading decisions and selects the offloading target if it decides to dispatch the task.

¹The code for reproducing the simulation results of this letter is available at: <https://github.com/ImanRHT/QECO>.

TABLE II
SIMULATION PARAMETERS

Parameter	Value
Computation capacity of MD f_i	2.6 GHz
Computation capacity of EN f_j^E	42.8 GHz
Transmission capacity of MD $r_{i,j}(t)$	14 Mbps
Task arrival rate	150 Task/sec
Size of task $\lambda_i(t)$	$\{1.0, 1.1, \dots, 7.0\}$ Mbits
Required CPU cycles of task $\rho_i(t)$	$\{0.197, 0.297, 0.397\} \times 10^3$
Deadline of task Δ_i	10 time slots (1 Sec)
Energy level percentage of MD $\phi_i(t)$	$\{25, 50, 75\}$
Computation power of EN p_j^E	5 W
Transmission power of MD p_i^T	2.3 W
Standby power of MD p_i^I	0.1 W

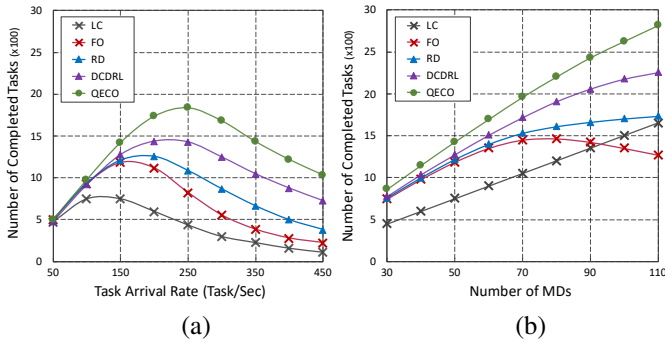


Fig. 3. The number of completed tasks under different computation workloads: (a) task arrival rate; (b) the number of MDs.

- 4) **DCDRL [14]**: This method is a distributed optimization algorithm designed for delay-sensitive tasks in a queuing-based MEC environment, where MDs strategically interact with multiple ENs. We select DCDRL as a benchmark due to its similarity to our work. It is based on the actor-critic framework [39], which underpins many state-of-the-art DRL algorithms, such as DDPG [40], and Proximal Policy Optimization [41].

B. Performance Comparison and Convergence

We first evaluate the number of completed tasks when comparing our proposed QECO algorithm with the other four schemes. As illustrated in Fig. 3 (a), the QECO algorithm

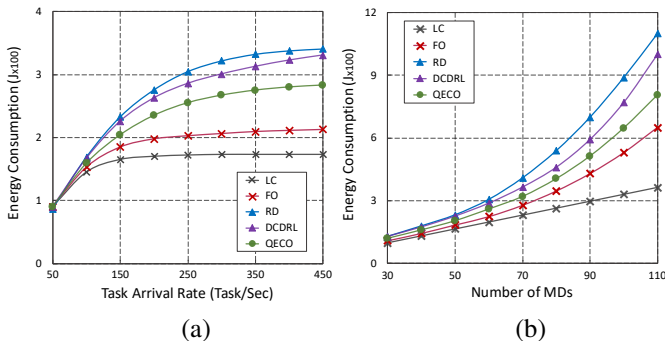


Fig. 4. The overall energy consumption under different computation workloads: (a) task arrival rate; (b) the number of MDs.

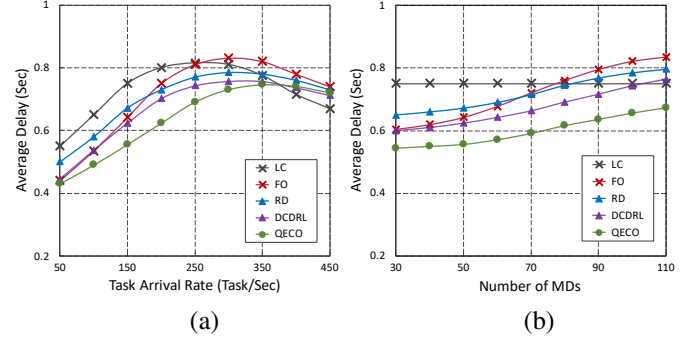


Fig. 5. The average delay under different computation workloads: (a) task arrival rate; (b) the number of MDs.

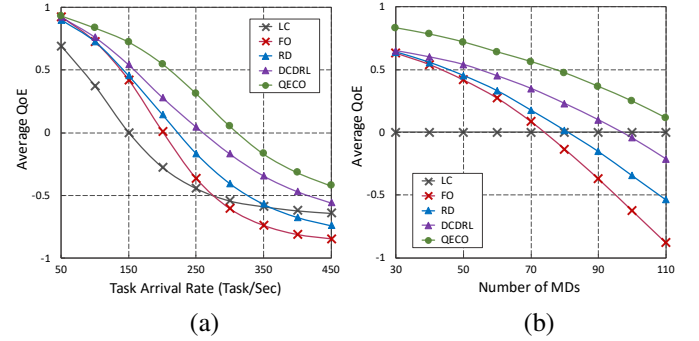


Fig. 6. The average QoE under different computation workloads: (a) task arrival rate; (b) the number of MDs.

consistently outperforms the benchmark methods when we vary the task arrival rate. At a lower task arrival rate (i.e., 50), most of the methods demonstrate similar proficiency in completing tasks. However, as the task arrival rate increases, the efficiency of QECO becomes more evident. Specifically, when the task arrival rate increases to 250, our algorithm can increase the number of completed tasks by 69.8% and 29.6% compared to RD and DCDRL, respectively. Similarly, in Fig. 3 (b), as the number of MDs increases, QECO shows significant improvements in the number of completed tasks compared to other methods, especially when faced with a large number of MDs. When there are 110 MDs, our proposed algorithm can effectively increase the number of completed tasks by at least 24.8% comparing with other methods. This achievement is attributed to the QECO's ability to effectively handle unknown workloads and prevent congestion at the ENs.

Figs. 4 (a) and 4 (b) illustrate the overall energy consumption for different values of task arrival rate and the number of MDs, respectively. At the lower task arrival rate, the total energy consumption of all methods is close to each other. The total energy consumption increases when we have a higher task arrival rate. As can be observed from Fig. 4 (a), at task arrival rate 450, QECO effectively reduces overall energy consumption by 18.6% and 16.2% compared to RD and DCDRL, respectively, as it takes into account the **energy level** of the MD in its decision-making process. However, it consumes more energy compared to LC and FO because they do not utilize all computing resources. In particular, LC only uses the MD's computational resources, while FO utilizes the allocated EN computing resources. —

In Fig. 4 (b), an increasing trend in overall energy consumption is observed as the number of MDs increases since the number of resources available in the system increases, which leads to higher energy consumption. The QECO algorithm consistently outperforms RD and DCDRL methods in overall energy consumption, especially when there are a large number of MDs. Specifically, QECO demonstrates a 27.4% and 20.5% reduction in overall energy consumption compared to RD and DCDRL, respectively, when the number of MDs increases to 110.

As shown in Fig. 5 (a), the QECO algorithm maintains a lower average delay compared to other methods as the task arrival rate increases from 50 to 350. Specifically, when the task arrival rate is 200, it reduces the average delay by at least 12.4% compared to other methods. However, for task arrival rates exceeding 350, QECO may experience a higher average delay compared to some of the other methods. This can be attributed to the fact that the other algorithms drop more tasks while our proposed algorithm is capable of completing a higher number of tasks, potentially leading to an increase in average delay. In Fig. 5 (b), as the number of MDs increases, we observe a rising trend in the average delay. It can be inferred that an increase in computational load in the system can lead to higher queuing delays and computations at ENs. Considering the QECO's ability to schedule workloads, when the number of MDs increases from 30 to 110, it consistently maintains a lower average delay which is at least 10.8% less than the other methods.

We further investigate the overall improvement achieved by the QECO algorithm in comparison to other methods in terms of the average QoE. This metric signifies the advantages MDs obtain by utilizing different algorithms. Fig. 6 (a) shows the average QoE for different values of the task arrival rate. This figure highlights the superiority of the QECO algorithm in providing MDs with an enhanced experience. At lower task arrival rates (i.e., 50-150), QECO performs optimally and maintains an average QoE of at least 0.72, while the other methods experience a steeper decline, with average QoE dropping to 0-0.54. Specifically, at a task arrival rate of 200, QECO improves the average QoE by 91.7% compared to DCDRL in the positive average QoE range. As task arrival rates increase to 300, the average QoE significantly decreases for all methods due to increased competition for resources in the MEC system. However, QECO still maintains a positive average QoE, while other methods fall to negative values, with a decline of up to -0.17. For higher task arrival rates (i.e., 350-450), QECO experiences a negative average QoE, but it shows at least a 42.6% and 24.6% smaller reduction in average QoE compared to RD and DCDRL, respectively.

Fig. 6 (b) illustrates the average QoE when we increase the number of MDs. The EN's workload grows when there are a larger number of MDs, leading to a reduction in the average QoE of all methods except LC. However, QECO effectively manages the uncertain load at the ENs. When the number of MDs increases from 30 to 110, QECO consistently maintains at least a 29.8% higher QoE compared to the other methods. Specifically, at a moderate number of 70 MDs, QECO achieves an average QoE of 0.57, showing a

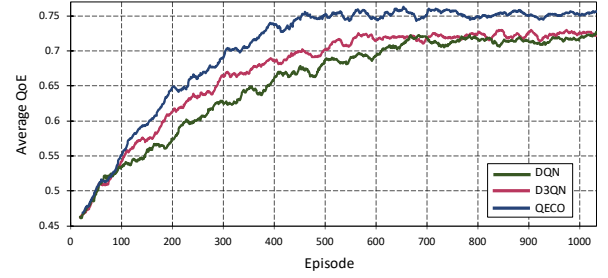


Fig. 7. The convergence of the average QoE across episodes under different DQN-based methods.

60.3% improvement compared to DCDRL. It is worth noting that although improvements in each of the QoE factors can contribute to enhancing system performance, it is essential to consider the user's demands in each time slot. Therefore, the key difference between QECO and other methods is that it prioritizes users' demands, enabling it to strike an appropriate balance among them, ultimately leading to a higher QoE for MDs.

We finally delve into the investigation of the convergence performance of the QECO algorithm in Fig. 7 and Fig. 8. To validate the effectiveness of the QECO algorithm, we assess its convergence speed and rate compared to the vanilla DQN and D3QN-based computation offloading methods, measured by the average QoE across episodes. As shown in Fig. 7, the average QoE for all computation offloading algorithms rises consistently as the episodes progress, until reaching convergence. Specifically, the vanilla DQN, D3QN, and our proposed algorithm converge after approximately 650, 550, and 400 iterations, respectively. This demonstrates the relatively faster convergence of the QECO algorithm, which also achieves a higher average QoE than the other methods. These results highlight the positive impact of workload prediction by the LSTM network and its advantage in efficiently utilizing the processing capabilities of MDs and ENs.

Furthermore, we explore the impact of two main hyperparameters on the convergence speed and the converged result of the proposed algorithm. Fig. 8 (a) illustrates the convergence of the proposed algorithm under different learning rates, where the learning rate regulates the step size per iteration towards minimizing the loss function. The QECO algorithm achieves an average QoE of 0.75 when the learning rate is 0.001, indicating relatively rapid convergence. However, with smaller learning rates (e.g., 0.0001) or larger values (e.g., 0.01), a slower convergence is observed. Fig. 8 (b) shows the convergence of the proposed algorithm under different batch sizes, which refer to the number of sampled experiences in each training round. An improvement in convergence performance is observed as the batch size increases from 4 to 16. However, further increasing the batch size from 16 to 32 does not notably enhance the converged QoE or convergence speed. Hence, a batch size of 16 may be more appropriate for training processes.

VII. CONCLUSION

In this paper, we focused on addressing the challenge of offloading in MEC systems, where strict task processing deadlines and energy constraints adversely impact system performance.

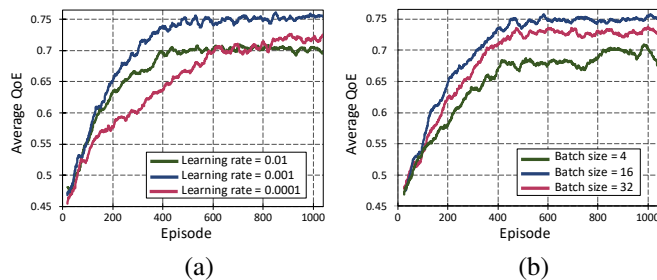


Fig. 8. The convergence of the average QoE across episodes under different hyper-parameters: (a) Learning rate; (b) Batch size.

We formulated an optimization problem that aims to maximize the QoE of each MD individually, while QoE reflects the energy consumption and task completion delay. To address the dynamic and uncertain mobile environment, we proposed a QoE-oriented DRL-based computation offloading algorithm called QECO. Our proposed algorithm empowers MDs to make offloading decisions without relying on knowledge about task models or other MDs' offloading decisions. The QECO algorithm not only adapts to the uncertain dynamics of load levels at ENs, but also effectively manages the ever-changing system environment. Through extensive simulations, we showed that QECO outperforms several established benchmark techniques, while demonstrating a rapid training convergence. Specifically, QECO increases the average user's QoE by 37% compared to an existing work. This advantage can lead to improvements in key performance metrics, including task completion rate, task delay, and energy consumption, under different system conditions and varying user demands.

There are multiple directions for future work. A complementary approach involves extending the task model by considering interdependencies among tasks. This can be achieved by incorporating a task call graph representation. Furthermore, in order to accelerate the learning of optimal offloading policies, it will be beneficial to take advantages of federated learning techniques in the training process. This will allow MDs to collectively contribute to improving the offloading model and enable continuous learning when new MDs join the network.

REFERENCES

- [1] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," *IEEE Commun. Surv. Tutor.*, vol. 19, no. 4, pp. 2322–2358, Aug 2017.
- [2] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge intelligence: Paving the last mile of artificial intelligence with edge computing," *Proc IEEE*, vol. 107, no. 8, pp. 1738–1762, Aug 2019.
- [3] A. Yousefpour, C. Fung, T. Nguyen, K. Kadiyala, F. Jalali, A. Niakanlahiji, J. Kong, and J. P. Jue, "All one needs to know about fog computing and related edge computing paradigms: A complete survey," *J. Syst. Archit.*, vol. 98, pp. 289–330, Sep 2019.
- [4] A. Kaur and A. Godara, "Machine learning empowered green task offloading for mobile edge computing in 5G networks," *IEEE Trans. Netw. Sci. Eng.*, vol. 21, no. 1, pp. 810–820, Feb 2024.
- [5] H. Shah-Mansouri and V. W. Wong, "Hierarchical fog-cloud computing for IoT systems: A computation offloading game," *IEEE Internet Things J.*, vol. 5, no. 4, pp. 3246–3257, May 2018.
- [6] C. Jiang, X. Cheng, H. Gao, X. Zhou, and J. Wan, "Toward computation offloading in edge computing: A survey," *IEEE Access*, vol. 7, pp. 131 543–131 558, Aug 2019.
- [7] L. Wu, P. Sun, H. Chen, Y. Zuo, Y. Zhou, and Y. Yang, "NOMA-enabled multiuser offloading in multicell edge computing networks: A coalition game based approach," *IEEE Trans. Netw. Sci. Eng.*, vol. 11, no. 2, pp. 2170–2181, Mar 2024.
- [8] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb 2015.
- [9] M. Guo, Q. Li, Z. Peng, X. Liu, and D. Cui, "Energy harvesting computation offloading game towards minimizing delay for mobile edge computing," *Computer Networks*, vol. 204, p. 108678, 2022.
- [10] Q. Tang, L. Liu, C. Jin, J. Wang, Z. Liao, and Y. Luo, "An uav-assisted mobile edge computing offloading strategy for minimizing energy consumption," *Computer Networks*, vol. 207, p. 108857, 2022.
- [11] Y. Li, B. Yang, H. Wu, Q. Han, C. Chen, and X. Guan, "Joint offloading decision and resource allocation for vehicular fog-edge computing networks: A contract-stackelberg approach," *IEEE Internet Things J.*, vol. 9, no. 17, pp. 15 969–15 982, 2022.
- [12] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Dueling network architectures for deep reinforcement learning," in *Proc. of International Conference on Machine Learning*. New York, NY, Jun 2016.
- [13] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, Sep 1997.
- [14] X. Qiu, W. Zhang, W. Chen, and Z. Zheng, "Distributed and collective deep reinforcement learning for computation offloading: A practical perspective," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 5, pp. 1085–1101, 2020.
- [15] J. Zhang, F. Shen, L. Tang, F. Yan, F. Qin, and C. Wang, "A multi-agent reinforcement learning approach for dynamic offloading with partial information-sharing in iot networks," in *2023 IEEE 98th Vehicular Technology Conference (VTC2023-Fall)*. IEEE, 2023, pp. 1–5.
- [16] S.-Y. Lin, C.-M. Huang, and T.-Y. Wu, "Multi-access edge computing-based vehicle-vehicle-rsu data offloading over the multi-rsu-overlapped environment," *IEEE Open J. Intell. Transp. Syst.*, vol. 3, pp. 7–32, 2022.
- [17] Z. Wei, B. Li, R. Zhang, X. Cheng, and L. Yang, "Many-to-many task offloading in vehicular fog computing: A multi-agent deep reinforcement learning approach," *IEEE Trans. Mob. Comput.*, vol. 23, no. 3, pp. 2107–2122, 2023.
- [18] S. Yuan, B. Dong, H. Lv, H. Liu, H. Chen, C. Wu, S. Guo, Y. Ding, and J. Li, "Adaptive incentivize for cross-silo federated learning in iiot: A multi-agent reinforcement learning approach," *IEEE Internet Things J.*, 2023.
- [19] L. Huang, S. Bi, and Y.-J. A. Zhang, "Deep reinforcement learning for online computation offloading in wireless powered mobile-edge computing networks," *IEEE Trans. Mob. Comput.*, vol. 19, no. 11, pp. 2581–2593, Jul 2019.
- [20] Z. Liu, Y. Zhao, J. Song, C. Qiu, X. Chen, and X. Wang, "Learn to coordinate for computation offloading and resource allocation in edge computing: A rational-based distributed approach," *IEEE Trans. Netw. Sci. Eng.*, vol. 9, no. 5, pp. 3136–3151, Dec 2021.
- [21] G. Wu, H. Wang, H. Zhang, Y. Zhao, S. Yu, and S. Shen, "Computation offloading method using stochastic games for software-defined-network-based multiagent mobile edge computing," *IEEE Internet Things J.*, vol. 10, no. 20, pp. 17 620–17 634, May 2023.
- [22] H. Zhou, K. Jiang, X. Liu, X. Li, and V. C. Leung, "Deep reinforcement learning for energy-efficient computation offloading in mobile-edge computing," *IEEE Internet Things J.*, vol. 9, no. 2, pp. 1517–1530, Jun 2021.
- [23] J. Chen, H. Xing, Z. Xiao, L. Xu, and T. Tao, "A drl agent for jointly optimizing computation offloading and resource allocation in mec," *IEEE Internet of Things Journal*, vol. 8, no. 24, pp. 17 508–17 524, 2021.
- [24] G. Wu, X. Chen, Z. Gao, H. Zhang, S. Yu, and S. Shen, "Privacy-preserving offloading scheme in multi-access mobile edge computing based on madrl," *J. Parallel Distributed Comput.*, vol. 183, p. 104775, Jan 2024.
- [25] Y. Dai, K. Zhang, S. Maharjan, and Y. Zhang, "Edge intelligence for energy-efficient computation offloading and resource allocation in 5G beyond," *IEEE Trans. Veh. Technol.*, vol. 69, no. 10, pp. 12 175–12 186, Aug 2020.
- [26] N. Zhao, Y.-C. Liang, D. Niyato, Y. Pei, M. Wu, and Y. Jiang, "Deep reinforcement learning for user association and resource allocation in heterogeneous cellular networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5141–5152, Aug 2019.
- [27] L. Liao, Y. Lai, F. Yang, and W. Zeng, "Online computation offloading with double reinforcement learning algorithm in mobile edge computing," *J Parallel Distrib Comput*, vol. 171, pp. 28–39, Jan 2023.
- [28] C. Sun, X. Li, C. Wang, Q. He, X. Wang, and V. C. Leung, "Hierarchical deep reinforcement learning for joint service caching and computation

- offloading in mobile edge-cloud computing,” *accepted for publication in IEEE Trans. Services Computing*, 2024.
- [29] Y. Gong, H. Yao, J. Wang, M. Li, and S. Guo, “Edge intelligence-driven joint offloading and resource allocation for future 6G Industrial Internet of Things,” *accepted for publication in IEEE Trans. Netw. Sci. Eng.*, 2024.
 - [30] Z. Gao, L. Yang, and Y. Dai, “Large-scale computation offloading using a multi-agent reinforcement learning in heterogeneous multi-access edge computing,” *IEEE Trans. Mob. Comput.*, vol. 22, no. 6, pp. 3425–3443, Jan 2023.
 - [31] M. Tang and V. W. Wong, “Deep reinforcement learning for task offloading in mobile edge computing systems,” *IEEE Trans. Mob. Comput.*, vol. 21, no. 6, pp. 1985–1997, Nov 2020.
 - [32] G. Wu, X. Chen, Y. Shen, Z. Xu, H. Zhang, S. Shen, and S. Yu, “Combining lyapunov optimization with actor-critic networks for privacy-aware iiot computation offloading,” *IEEE Internet Things J.*, Jan 2024.
 - [33] G. Wu, Z. Xu, H. Zhang, S. Shen, and S. Yu, “Multi-agent drl for joint completion delay and energy consumption with queuing theory in mec-based iiot,” *J. Parallel Distributed Comput.*, vol. 176, pp. 80–94, Jun 2023.
 - [34] H. Huang, Q. Ye, and Y. Zhou, “Deadline-aware task offloading with partially-observable deep reinforcement learning for multi-access edge computing,” *IEEE Trans. Netw. Sci. Eng.*, vol. 9, no. 6, pp. 3870–3885, Sep 2021.
 - [35] M. Bolourian and H. Shah-Mansouri, “Deep Q-learning for minimum task drop in SWIPT-enabled mobile-edge computing,” *IEEE Wireless Commun. Letters*, vol. 13, no. 3, pp. 894–898, Mar 2024.
 - [36] M. S. Munir, N. H. Tran, W. Saad, and C. S. Hong, “Multi-agent meta-reinforcement learning for self-powered and sustainable edge computing systems,” *IEEE Trans. Netw. Serv. Manag.*, vol. 18, no. 3, pp. 3353–3374, 2021.
 - [37] Y. Mao, J. Zhang, and K. B. Letaief, “Dynamic computation offloading for mobile-edge computing with energy harvesting devices,” *IEEE J. Sel. Areas Commun.*, vol. 34, no. 12, pp. 3590–3605, Sep 2016.
 - [38] A. Parekh and R. G. Gallager, “A generalized processor sharing approach to flow control in integrated services networks: The single-node case,” *IEEE/ACM Trans. Netw.*, vol. 1, no. 3, pp. 344–357, Jun 1993.
 - [39] V. Konda and J. Tsitsiklis, “Actor-critic algorithms,” *Advances in neural information processing systems*, vol. 12, 1999.
 - [40] T. Lillicrap, “Continuous control with deep reinforcement learning,” *arXiv preprint arXiv:1509.02971*, 2015.
 - [41] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” *arXiv preprint arXiv:1707.06347*, 2017.