

# QOCO: 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 responsive and reliable services. This challenge stands as one of the most primary key factors contributing to handling dynamic load 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 decentralized QoE-oriented computation offloading (QOCO) 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 QOCO. Simulation results validate that the QOCO algorithm efficiently exploits the computational resources of edge nodes. Consequently, it can complete 14% more tasks and reduce task delay and energy consumption by 9% and 6%, respectively. These together contribute to a significant improvement of at least 37% in average QoE compared to an existing algorithm.

**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 battery 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]. 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 [5]. In particular, determining the optimal offloading strategy, scheduling the tasks, and selecting the most suitable EN for task offloading are main the challenges that demand careful consideration.

Researchers have recently proposed several task offloading algorithms to tackle the aforementioned issues. Mao *et al.* in [6] introduced a computation offloading algorithm for

MEC systems. This scheme aims to reduce the MD's energy consumption while meeting the computation delay constraints. In [7], Zhang *et al.* proposed an offloading scheme for MEC in heterogeneous cellular networks, taking into account the energy-constrained MEC scenario and the heterogeneity of MDs when optimizing the offloading decisions. Bi *et al.* in [8] proposed an algorithm to optimize decision-making about computation offloading and power transferring in a wireless-powered MEC. Ning *et al.* in [9] introduced a heuristic algorithm designed to enhance the efficiency of the partial computation offloading model. This algorithm effectively adapts offloading decisions dynamically. The works in [6]-[9] primarily focus on quasi-static systems and are not well-suited for dynamic systems and time-varying conditions. Furthermore, the uncertain requirements and sensitive latency properties of computation tasks in MEC systems with limited resources pose nontrivial challenges that can significantly impact computation offloading performance.

To cope with the dynamic nature of the network, recent research has proposed several DRL-based algorithms. DRL combines reinforcement learning principles with deep neural networks to determine optimal decision-making policies [10]. By capturing the dynamics of environments and learning strategies for accomplishing long-term objectives through repeated interactions, DRL can effectively tackle the challenges arising from the ever-changing nature of networks, MDs, and servers' heterogeneity. This ultimately improves the MD users' QoE. In [11], Huang *et al.* focused on a wireless-powered MEC network. They proposed a DRL-based online approach, capable of attaining near-optimal decisions. This is achieved by selectively considering a compact subset of candidate actions in each iteration. Liao *et al.* in [12] introduced a double reinforcement learning algorithm for performing online computation offloading in MEC. This algorithm optimizes transmission power and scheduling of CPU frequency to minimize both task computation delay and energy consumption. In [13], Zhao *et al.* proposed a computation offloading algorithm based on DRL, which addresses the competition for wireless channels to optimize long-term downlink utility. In this approach, each MD requires quality-of-service information from other MDs. Tang *et al.* in [14] 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 [15], Chen

*et al.* introduced an online algorithm with polynomial time complexity for offloading in MEC. This algorithm achieves a close-to-optimal energy consumption while ensuring that the delay is bounded. In [16], Dai *et al.* introduced the integration of action refinement into DRL and designed an algorithm to concurrently optimize resource allocation and computation offloading. Zhou *et al.* in [17] used an MDP to study MEC and modeled the interactions between the policy-system environment. They proposed a Q-learning approach to achieve optimal resource allocation strategies and computation offloading. In [18], Gao *et al.* introduced an attention-based multi-agent algorithm designed for decentralized computation offloading. This algorithm effectively tackles the challenges of dynamic resource allocation in large-scale heterogeneous networks.

Although DRL-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 as assumed in [13]-[14], but also by energy consumption. A more comprehensive approach is required to address the dynamic requirements of individual users in real-time scenarios with multiple MDs and ENs. Nevertheless, prior research in [11]-[18] fail to tackle some or all of these challenges.

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 performance. We propose the distributed QOCO algorithm that leverages DRL to efficiently handle task offloading in uncertain loads at ENs. This algorithm empowers MDs to make offloading decisions utilizing only locally observed information, such as task size, queue details, battery status, and historical workloads at the ENs. By adopting the appropriate policy based on each MD's specific requirements at any given time, the QOCO 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 as an MDP for time-sensitive tasks. This approach takes into account the dynamic nature of workloads at the ENs and concentrates on providing high performance in the MEC system while maximizing the long-term QoE.
- **DRL-based Offloading Algorithm:** To address the problem of long-term QoE maximization, we focus on task completion, task delay, and energy consumption to quantify the MDs' QoE. We propose QOCO 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 battery level, our approach leverages deep Q-network (DQN)[19] and long short-term memory(LSTM) [20] to prioritize and strike an appropriate balance between QoE factors.
- **Performance Evaluation:** We conduct comprehensive experiments to evaluate the QOCO's performance. The

simulation results confirm that when compared with the potential game-based offloading algorithm (PGOA) in [21] and several benchmark methods, our algorithm effectively utilizes the processing capabilities of MDs and ENs, resulting in a substantial improvement in the average QoE of the system.

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

## II. SYSTEM MODEL

We investigate an 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 time slots denoted by  $\mathcal{T} = \{1, 2, \dots, T\}$ , with  $T$  time slots, each representing a duration of  $\tau$  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 from a discrete set  $\Lambda = \{\lambda_1, \lambda_2, \dots, \lambda_\theta\}$ , where  $\theta$  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 queues 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 at time slot  $t \in \mathcal{T}$ , then it will either be dispatched or dropped at 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. The value of  $\delta_i^C(t)$  is calculated as follows:

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

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  $t + \Delta_i(t) - 1$ .

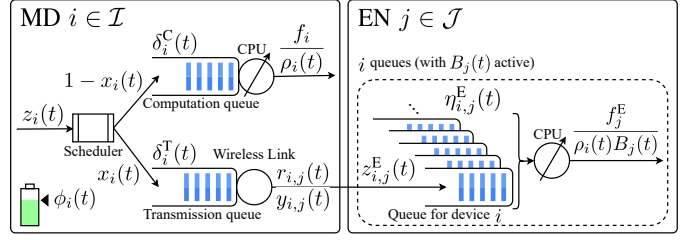


Fig. 1. An illustration of MD  $i \in \mathcal{I}$  and EN  $j \in \mathcal{J}$  in the MEC system.

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 [6].

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  $j$  at 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  $i^{\text{th}}$  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  $t$ , given by:

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

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 [22]. 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:

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

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.

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

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

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  $l_{i,j}^E(t)$  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

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

where  $p_i^I$  is the energy consumption of MD  $i \in \mathcal{I}$  in the standby state. Accordingly, by combining (4), (15) and (16), the overall energy consumption associated with offloading task  $z_i(t)$  is formulated as

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

### III. 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 battery level, which plays a crucial role in decision-making. Specifically, when an MD observes its state (e.g. task size, queue details, and battery level) 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 battery 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:

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

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 battery level of MD  $i$  in time slot  $t$ . Considering the power modes of a 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 a matrix with a size of  $T^s \times J$ . 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.

### C. QoE Function

The QoE function evaluates the influence of agent's actions by taking several key performance factors into account. 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:

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

where  $\mathcal{D}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)) = 0$  when task  $z_i(t)$  is dropped. Correspondingly, we denote the energy consumption of task  $z_i(t)$  when taking action  $\mathbf{a}_i(t)$  in the observed state  $\mathbf{s}_i(t)$  as  $\mathcal{E}_i(\mathbf{s}_i(t), \mathbf{a}_i(t))$ , which is:

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

Given the delay and energy consumption of task  $z_i(t)$ , we also define  $\mathcal{C}_i(\mathbf{s}_i(t), \mathbf{a}_i(t))$  that denotes the associate cost of task  $z_i(t)$  given the action  $\mathbf{a}_i(t)$  in the state  $\mathbf{s}_i(t)$ .

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

where  $\phi_i(t)$  represents the MD  $i$ 's battery level. When the MD is operating in performance mode, the primary focus is on minimizing task delays, thus the delay contributes more to the cost. On the other hand, when the MD switches to ultra power-saving mode, the main attention is directed toward reducing power consumption.

Finally, we define  $\mathbf{q}_i(\mathbf{s}_i(t), \mathbf{a}_i(t))$  as the QoE associated with task  $z_i(t)$  given the selected action  $\mathbf{a}_i(t)$  and the observed state  $\mathbf{s}_i(t)$ . The QoE function is defined as follows:

$$\begin{aligned} \mathbf{q}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)) &= \\ \begin{cases} \mathcal{R} - \mathcal{C}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)) & \text{if task } z_i(t) \text{ is processed,} \\ -\mathcal{E}_i(\mathbf{s}_i(t), \mathbf{a}_i(t)) & \text{if task } z_i(t) \text{ is dropped,} \end{cases} \end{aligned} \quad (24)$$

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

### 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, the MD  $i$  determines an action  $\mathbf{a}_i(t) \in \mathcal{A}$ , according to policy  $\pi_i$  given the observed environment state  $\mathbf{s}_i(t) \in \mathcal{S}$ . The MD  $i$  aims to find its optimal policy  $\pi_i^*$  which maximizes the long-term QoE,

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

where  $\gamma \in (0, 1]$  is a discount factor and determines the balance between instant QoE and long-term QoE. As  $\gamma$  approaches 0, the MD prioritizes QoE within the current time slot exclusively. Conversely, as  $\gamma$  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.

## IV. DRL-BASED OFFLOADING ALGORITHM

We now present QOCO 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 solution

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  $\theta_i$  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  $\boldsymbol{\eta}_i^E(t-1)$  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 [20]. This network maintains a memory state  $\mathcal{H}(t)$  that evolves over time, enabling 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 track the variations in sequences from  $\mathbf{h}_1(t)$  to  $\mathbf{h}_{T^s}(t)$ , where vector  $\mathbf{h}_i(t) = (h_{i,j}(t))_{j \in \mathcal{J}}$ , 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 [23] 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  $V_i(\mathbf{s}_i(t); \theta_i)$  and  $A_i(\mathbf{s}_i(t), \mathbf{a}; \theta_i)$  to denote the state-value and the action-advantage value of action  $\mathbf{a} \in \mathcal{A}$  under state  $\mathbf{s}_i(t) \in \mathcal{S}$ , respectively. The parameter  $\theta_i$  is responsible for determining these values, and it can be adjusted when training the QOCO 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  $\theta_i$  establishes a functional relationship that maps Q-values to pairs of state-action.

### B. QoE-Oriented DRL-Based Algorithm

The QOCO 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

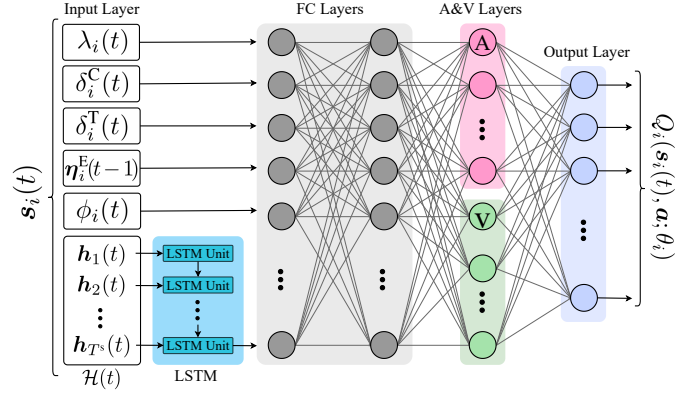


Fig. 2. The neural network of MD  $i \in \mathcal{I}$ , which characterize the Q-value of each action  $\mathbf{a} \in \mathcal{A}$  under state  $\mathbf{s}_i(t) \in \mathcal{S}$ .

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  $\mathcal{M}_i$  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  $\theta_i^E$  and  $\theta_i^T$ , 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  $\theta_i^E$ . 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  $N^{\text{ep}}$  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 parameter vector  $\theta_i^E$  of  $Net_i^E$  from EN  $j_i$ , MD  $i$  chooses the following action for task  $z_i(t)$ .

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



**Algorithm 1** QOCO Algorithm (Offloading Decision)**Input:** state space  $\mathcal{S}$ , action space  $\mathcal{A}$ **Output:** MD  $i \in \mathcal{I}$  experience

```

1: for episode 1 to  $N^{\text{ep}}$  do
2:   Initialize  $s_i(1)$ 
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  $a_i(t)$  based on (27);
8:     end if
9:     Observe a set of QoEs  $\{q_i(t'), t' \in \mathcal{F}_i^t\}$ ;
10:    Observe the next state  $s_i(t+1)$ ;
11:    for each task  $z_i(t')$  where  $t' \in \mathcal{F}_i^t$  do
12:      Send  $(s_i(t'), a_i(t'), q_i(t'), s_i(t'+1))$  to EN  $j_i$ ;
13:    end for
14:  end for
15: end for

```

where w.p. stands for with probability, and  $\epsilon$  represents the random exploration probability. The value of  $Q_i^E(s_i(t), a; \theta_i^E)$  indicates the Q-value under the parameter  $\theta_i^E$  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  $s_i(t)$ .

In the next time slot  $t+1$ , MD  $i$  observes the state  $s_i(t+1)$ . However, due to the potential for tasks to extend across multiple time slots, QoE  $q_i(t)$  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:

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

Therefore, MD  $i$  observes a set of QoEs  $\{q_i(t') \mid t' \in \mathcal{F}_i^t\}$  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  $(s_i(t), a_i(t), q_i(t), s_i(t+1))$  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  $\mathcal{M}_i$  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  $\theta_i^E$ , obtained from  $Net_i^E$ , back to MD  $i$ . On the other side, if EN  $j$  receives an experience  $(s_i(t), a_i(t), q_i(t), s_i(t+1))$  from MD  $i \in \mathcal{I}_j$ , the EN stores this experience in the replay buffer  $\mathcal{M}_i$  associated with that MD.

The EN randomly selects a sample collection of experiences from the replay buffer, denoted as  $\mathcal{N}$ . For each experience

**Algorithm 2** QOCO Algorithm (Training Process)

```

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 from MD  $i \in \mathcal{I}_j$  then
9:     Store  $(s_i(t'), a_i(t'), q_i(t'), s_i(t'+1))$  in  $\mathcal{M}_i$ ;
10:    Get a collection of experiences  $\mathcal{I}$  from  $\mathcal{M}_i$ ;
11:    for each experience  $i \in \mathcal{I}$  do
12:      Get experience  $(s_i(n), a_i(n), q_i(n), s_i(n+1))$ ;
13:      Generate  $\hat{Q}_{i,n}^T$  according to (28);
14:    end for
15:    Set vector  $\hat{\mathbf{Q}}_i^T := (\hat{Q}_{i,n}^T)_{n \in \mathcal{N}}$ ;
16:    Update  $\theta_i^E$  to minimize  $L(\theta_i^E, \hat{\mathbf{Q}}_i^T)$  in (30);
17:    Count := Count + 1;
18:    if mod(Count, ReplaceThreshold) = 0 then
19:       $\theta_i^T := \theta_i^E$ ;
20:    end if
21:  end if
22: end while
23: end while

```

$n \in \mathcal{N}$ , 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 set  $\mathcal{N}$ , the EN trains the MD's neural network using previous sample experiences. Simultaneously, it updates  $\theta_i^E$  in  $Net_i^E$  and computes vector  $\hat{\mathbf{Q}}_i^T = (\hat{Q}_{i,n}^T)_{n \in \mathcal{N}}$ . 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  $\theta_i^T$  in  $Net_i^T$ , which enhances the approximation of the long-term QoE when computing the target Q-values in (28).

TABLE I  
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 $\Delta_i$	10 time slots (1 Sec)
Battery 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

## V. PERFORMANCE EVALUATION

In this section, we first present the simulation setup and training configuration. We then evaluate the performance of QOCO in comparison to three baseline schemes in addition to the existing work [21].

### A. Simulation Setup

We consider a computation task offloading scenario involving 50 MDs and 5 ENs in the MEC environment, similar to [14], and we follow the model presented in [17] to determine the energy consumption. All the parameters are given in Table I. To train the MDs' neural networks, we adopt a scenario comprising 1000 episodes. Each episode contains 100 time slots, equivalent to 0.1 second. QOCO algorithm performs online training, incorporating real-time experience to continually improve the offloading strategy. Specifically, we employ a batch size of 16, maintain a fixed learning rate of 0.001, and set the discount factor  $\gamma$  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.

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 to ENs and selects 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.
- 4) **PGOA [21]:** This existing method is a distributed optimization algorithm designed for delay-sensitive tasks in an environment where MDs interact strategically with multiple ENs. We select PGOA as a benchmark method due to its similarity to our work.

### B. Performance Comparison

We first evaluate the number of completed tasks when comparing our proposed QOCO algorithm with the other

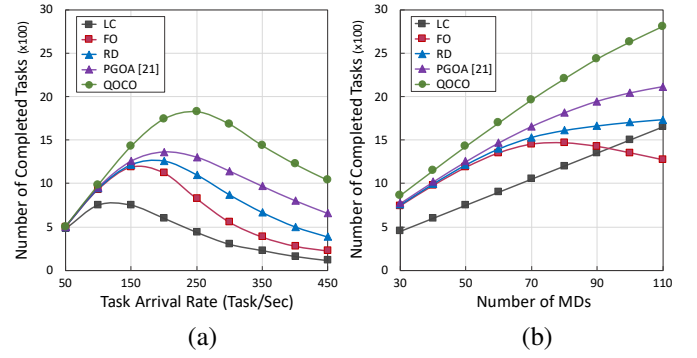


Fig. 3. The number of completed tasks under different computation workloads: (a) task arrival rate; (b) the number of MDs.

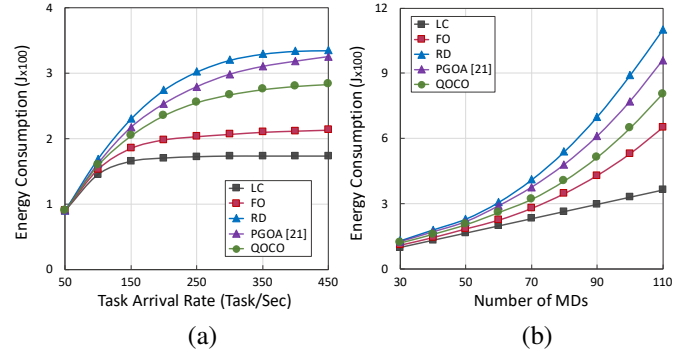


Fig. 4. The overall energy consumption under different computation workloads: (a) task arrival rate; (b) the number of MDs.

four schemes. In Fig. 3 (a), the QOCO algorithm 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 equal proficiency in completing tasks. However, as the task arrival rate increases, the efficiency of QOCO becomes more evident. Specifically, when the task arrival rate increases to 250, our algorithm can increase the number of completed tasks by 73% and 47% compared to RD and PGOA, respectively. Similarly, in Fig. 3 (b), as the number of MDs increases, QOCO shows significant improvements in the number of completed tasks compared to other methods, especially when faced with a large number of MDs. When the number of MDs reaches 110, our proposed algorithm can effectively increase the number of completed tasks by at least 34% comparing with other methods. This achievement is attributed to the QOCO's ability to effectively handle unknown workloads and prevent congestion at the ENs.

Fig. 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. When the task arrival rate increases, the total energy consumption of each method increases differently. It can be found that the production of more computational tasks leads to a continuous increase in energy consumption until the capacity is saturated, and the energy consumption remains constant at its maximum. The QOCO algorithm takes into account the battery level of the MD



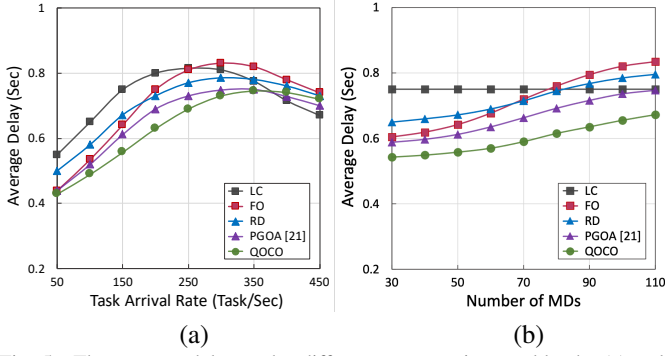


Fig. 5. The average delay under different computation workloads: (a) task arrival rate; (b) the number of MDs.

in its decision-making process. In Fig. 4 (a), at task arrival rate 450, QOCO effectively reduces overall energy consumption by 18% and 15% compared to RD and PGOA, respectively. However, it consumes more energy compared to LC and FC because they do not utilize the full capacity (LC uses the MD's computational resources, while FC utilizes the transmission link and allocated capacity from EN to MD). In Fig. 4 (b), an increasing trend in overall energy consumption is observed as the number of MDs increases. With the increase in MDs, the number of resources available in the system increases, which leads to higher energy consumption in the system. The QOCO algorithm consistently outperforms the benchmark methods in overall energy consumption, especially when the number of MDs is large. Specifically, QOCO demonstrates a 27% and 16% reduction in overall energy consumption compared to RD and PGOA, respectively, when the number of MDs increases to 110.

As shown in Fig. 5 (a), the QOCO 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% compared to other methods. However, for task arrival rates exceeding 350, QOCO may experience a higher average delay compared to some of the other methods. This can be attributed to the fact that our algorithm processes a larger number of completed tasks in the same time slot, 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 QOCO's ability to schedule workloads, when the number of MDs increases from 30 to 110, it consistently maintains an average delay at least 8% lower than the other methods.

We finally delve into the investigation of the overall improvement achieved by the QOCO 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 task arrival rate. This figure indicates the superiority of the QOCO algorithm in providing MDs with an enhanced experience.

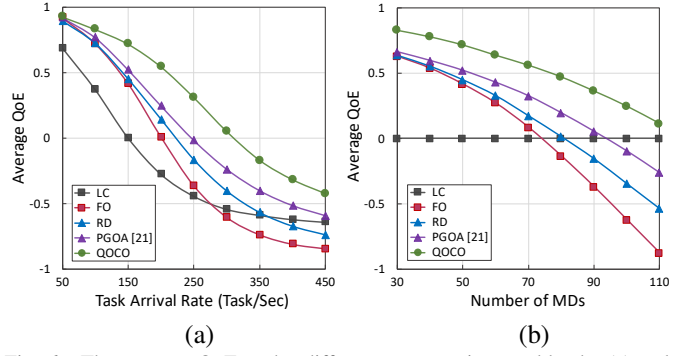


Fig. 6. The average QoE under different computation workloads: (a) task arrival rate; (b) the number of MDs.

Specifically, when the task arrival rate is moderate (i.e., 250), QOCO improves the average QoE by 57% and 33% compared to RD and PGOA, respectively. Fig. 6 (b) illustrates the average QoE when we increase the number of MDs. The higher the number of MDs, the workload on the ENs rises, leading to a reduction in the average QoE of all methods except LC. However, QOCO effectively manages the uncertain load at the ENs. When the number of MDs increases to 90, QOCO achieves at least a 29% higher QoE comparing with the other methods.

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 QOCO 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.

## VI. 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. 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 QOCO. Our proposed algorithm empowers MDs to make offloading decisions without relying on knowledge about task models or other MDs' offloading decisions. The QOCO 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 demonstrated that QOCO outperforms several established benchmark techniques. Specifically, QOCO 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 call graph representation to develop dependency among task partitions. Furthermore, in order to accelerate the learning of optimal offloading policies, it will be beneficial to enable MDs to take advantage of federated learning techniques in the training process. This allows 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] 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.
- [5] 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.
- [6] 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.
- [7] K. Zhang, Y. Mao, S. Leng, Q. Zhao, L. Li, X. Peng, L. Pan, S. Maharjan, and Y. Zhang, "Energy-efficient offloading for mobile edge computing in 5G heterogeneous networks," *IEEE Access*, vol. 4, pp. 5896–5907, Aug 2016.
- [8] S. Bi and Y. J. Zhang, "Computation rate maximization for wireless powered mobile-edge computing with binary computation offloading," *IEEE Trans. Wirel. Commun.*, vol. 17, no. 6, pp. 4177–4190, Apr 2018.
- [9] Z. Ning, P. Dong, X. Kong, and F. Xia, "A cooperative partial computation offloading scheme for mobile edge computing enabled Internet of Things," *IEEE Internet Things J.*, vol. 6, no. 3, pp. 4804–4814, Sep 2018.
- [10] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep reinforcement learning: A brief survey," *IEEE Signal Process Mag.*, vol. 34, no. 6, pp. 26–38, Nov 2017.
- [11] 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.
- [12] 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.
- [13] 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. Wirel. Commun.*, vol. 18, no. 11, pp. 5141–5152, Aug 2019.
- [14] 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.
- [15] Y. Chen, N. Zhang, Y. Zhang, X. Chen, W. Wu, and X. Shen, "Energy efficient dynamic offloading in mobile edge computing for Internet of Things," *IEEE Trans. on Cloud Comput.*, vol. 9, no. 3, pp. 1050–1060, Feb 2019.
- [16] 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.
- [17] 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.
- [18] 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.
- [19] 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.
- [20] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput*, vol. 9, no. 8, pp. 1735–1780, Sep 1997.
- [21] L. Yang, H. Zhang, X. Li, H. Ji, and V. C. Leung, "A distributed computation offloading strategy in small-cell networks integrated with mobile edge computing," *IEEE ACM Trans Netw.*, vol. 26, no. 6, pp. 2762–2773, 2018.
- [22] A. K. 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.
- [23] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Dueling network architectures for deep reinforcement learning," in *International conference on machine learning*. PMLR, Jun 2016, pp. 1995–2003.