

Deep Reinforcement Learning for Task Offloading in Mobile Edge Computing Systems

Ming Tang and Vincent W.S. Wong

Abstract—In mobile edge computing systems, an edge node may have a high load when a large number of mobile devices offload their tasks to it. Those offloaded tasks may experience large processing delay or even be dropped when their deadlines expire. Due to the uncertain load dynamics at the edge nodes, it is challenging for each device to determine its offloading decision (i.e., whether to offload or not, and which edge node it should offload its task to) in a decentralized manner. In this work, we consider non-divisible and delay-sensitive tasks as well as edge load dynamics, and formulate a task offloading problem to minimize the expected long-term cost. We propose a model-free deep reinforcement learning-based distributed algorithm, where each device can determine its offloading decision without knowing the task models and offloading decision of other devices. To improve the estimation of the long-term cost in the algorithm, we incorporate the long short-term memory (LSTM), dueling deep Q-network (DQN), and double-DQN techniques. Simulation results with 50 mobile devices and five edge nodes show that the proposed algorithm can reduce the ratio of dropped tasks and average task delay by 86.4% – 95.4% and 18.0% – 30.1%, respectively, when compared with several existing algorithms.

Index Terms—Mobile edge computing, fog computing, computation offloading, resource allocation, deep reinforcement learning, deep Q-learning.



1 INTRODUCTION

1.1 Background and Motivation

Nowadays, mobile devices are responsible for processing more and more computational intensive tasks, such as data processing, artificial intelligence, and virtual reality. Despite the development of mobile devices, these devices may not be able to process all their tasks locally with a low latency due to their limited computational resources. To facilitate efficient task processing, mobile edge computing (MEC) [1], also known as fog computing [2] and multi-access edge computing [3], is introduced. MEC facilitates mobile devices to offload their computational intensive tasks to nearby edge nodes for processing in order to reduce the task processing delay. It can also reduce the ratio of dropped tasks for those delay-sensitive tasks.

In MEC, there are two main questions related to task offloading. The first question is whether a mobile device should offload its task to an edge node or not. The second question is that if a mobile device decides to perform offloading, then which edge node should the device offload its task to. To address these questions, some existing works have proposed task offloading algorithms. Wang *et al.* in [4] proposed an algorithm to determine the offloading decisions of the mobile devices to maximize the network revenue. Bi *et al.* in [5] focused on a wireless-powered MEC scenario and proposed an algorithm to jointly optimize the offloading and power transfer decisions. In these works [4], [5], the processing capacity that each mobile device obtained from an edge node is independent of the number of tasks offloaded to the edge node.

In practice, however, edge nodes may have limited processing capacities, so the processing capacity that an edge node allocated to a mobile device depends on the *load level* at the edge node (i.e., number of concurrent tasks offloaded to the edge node). When a large number of mobile devices offload their tasks to the same edge node, the load at that edge node can be high, and hence those offloaded tasks may experience large processing delay. Some of the tasks may even be dropped when their deadlines expire. Some existing works have addressed the load levels at the edge nodes and proposed centralized task offloading algorithms. Eshraghi *et al.* in [6] considered the uncertain computational requirements of the mobile devices, and proposed an algorithm that optimizes the offloading decisions of the mobile devices and the computational resource allocation decision of the edge node. Lyu *et al.* in [7] focused on delay-sensitive tasks and proposed an algorithm to minimize the task offloading energy consumption subject to the task deadline constraint. In [8], Chen *et al.* considered a software-defined ultra-dense network, and designed a centralized algorithm to minimize the task processing delay. In [9], Poularakis *et al.* studied the joint optimization of task offloading and routing, taking into account the asymmetric requirements of the tasks. These centralized algorithms in [6]–[9], however, may require global information of the system (e.g., the arrivals and the sizes of the tasks of all mobile devices) and may incur high signaling overhead.

Other works have proposed distributed task offloading algorithms considering the load levels at the edge nodes, where each mobile device makes its offloading decision in a decentralized manner. Note that designing such a distributed algorithm is challenging. This is because when a device makes an offloading decision, the device does not know *a priori* the load levels at the edge nodes, since the load also depends on the offloading decisions and task

Ming Tang and Vincent W.S. Wong are with the Department of Electrical and Computer Engineering, The University of British Columbia, Vancouver, Canada.

E-mail: {mingt, vincentw}@ece.ubc.ca

models (e.g., the size and arrival time of the task) of other mobile devices. In addition, the load levels at the edge nodes may change over time. To address these challenges, Lyu *et al.* in [10] focused on divisible tasks and proposed a Lyapunov-based algorithm to ensure the stability of the task queues. In [11], Li *et al.* considered the strategic offloading interaction among mobile devices and proposed a price-based distributed algorithm. Shah-Mansouri *et al.* in [12] designed a potential game-based offloading algorithm to maximize the quality-of-experience of each device. Jošilo *et al.* in [13] designed a distributed algorithm based on a Stackelberg game. Yang *et al.* in [14] proposed a distributed offloading algorithm to address the wireless channel competition among mobile devices. Neto *et al.* in [15] proposed an estimation-based method, where each device makes its offloading decision based on the estimated processing and transmission capacities.

In this work, we focus on the task offloading problem in an MEC system and propose a distributed algorithm that addresses the unknown load level dynamics at the edge nodes. Comparing with the aforementioned works [10]–[15], we consider a different and more realistic MEC scenario. First, the existing work [10] considered divisible tasks (i.e., tasks can be arbitrarily divided), which may not be realistic due to the dependency among the bits in a task. On the other hand, although the works [11]–[15] considered non-divisible tasks, they do not take into account the underlying queuing systems. As a result, the processing and transmission of each task should always be accomplished within one time slot (or before the arrival of the next task), which may not always be guaranteed in practice. Different from those works [10]–[15], we consider non-divisible tasks together with queuing systems and take into account the practical scenario where the processing and transmission of each task can continue for multiple time slots. This scenario is challenging to deal with, because for the task of a mobile device arrived in a time slot, its delay can be affected by the decisions of the tasks of other devices arrived in the previous time slots. Second, different from the related works [10]–[15] which considered delay-tolerant tasks, we take into account delay-sensitive tasks with processing deadlines. This is challenging to address, because the processing deadlines will affect the load level dynamics at the edge nodes and hence affect the delay of the offloaded tasks.

The model-free deep reinforcement learning (DRL) techniques (e.g., deep Q-learning [16]) are candidate methods for solving the task offloading problem in the MEC system, as these methods enable agents to make decisions based on local observations without estimating the dynamics involved in the model. Some existing works such as [17]–[19] have proposed DRL-based algorithms for the MEC system, while they focused on centralized offloading algorithms. Zhao *et al.* in [20] proposed a DRL-based distributed offloading algorithm that addresses the wireless channel competition among mobile devices, while the algorithm at each mobile device requires the quality-of-service information of other mobile devices. Different from those works [17]–[20], we aim to propose a DRL-based distributed algorithm that can address the unknown load dynamics at the edge nodes. It should also enable each mobile device to make its offloading decision without knowing the information (e.g., task mod-

els, offloading decisions) of other mobile devices.

1.2 Solution Approach and Contributions

In this work, we take into account the unknown load level dynamics at the edge nodes and propose a DRL-based distributed offloading algorithm for the MEC system. In the proposed algorithm, each mobile device can determine the offloading decision in a decentralized manner using the information observed locally, including the size of its task, the information of its queues, and the historical load levels at the edge nodes. In addition, the proposed algorithm can handle the time-varying system environments, including the arrival of new tasks, the computational requirement of each task, and the offloading decisions of other mobile devices.

The main contributions are as follows.

- *Task Offloading Problem for the MEC System:* We formulate a task offloading problem taking into account the load level dynamics at the edge nodes to minimize the expected long-term cost (considering the delay of the tasks and the penalties for those tasks being dropped). In this problem, we consider non-divisible and delay-sensitive tasks and use queuing systems to model the processing and transmission processes of the tasks.
- *DRL-based Task Offloading Algorithm:* To achieve the expected long-term cost minimization considering the unknown load dynamics at the edge nodes, we propose a model-free DRL-based distributed offloading algorithm that enables each mobile device to make its offloading decision without knowing the task models and offloading decisions of other mobile devices. To improve the estimation of the expected long-term cost in the proposed algorithm, we incorporate the long short-term memory (LSTM), dueling deep Q-network (DQN), and double-DQN techniques.
- *Performance Evaluation:* We perform simulations and show that when compared with the potential game based offloading algorithm (PGOA) in [14] and the user-level online offloading framework (ULOOF) in [15], our proposed DRL-based algorithm can better exploit the processing capacities of the mobile devices and edge nodes, and it can significantly reduce the ratio of dropped tasks and the average delay. Under a scenario with 50 mobile devices and five edge nodes, our proposed algorithm can reduce the ratio of dropped tasks by 86.4% – 95.4% and reduce the average delay by 18.0% – 30.1% when compared with the existing algorithms.

The rest of this paper is organized as follows. The system model is presented in Section 2, and the problem formulation is given in Section 3. We present the DRL-based algorithm in Section 4 and evaluate its performance in Section 5. Conclusions are given in Section 6. For notation, we use \mathbb{Z}_{++} to denote the set of positive integers.

2 SYSTEM MODEL

We consider a set of edge nodes $\mathcal{N} = \{1, 2, \dots, N\}$ and a set of mobile devices $\mathcal{M} = \{1, 2, \dots, M\}$ in an MEC system.

The mobile devices can offload their computational tasks to the edge nodes for processing. We consider one episode that contains a set of time slots $\mathcal{T} = \{1, \dots, T\}$, where each time slot has a duration of Δ seconds. In the following, we present the mobile device model and the edge node model, respectively, with an illustration given in Fig. 1.

2.1 Mobile Device Model

During each time slot, we assume that a mobile device either has a new task arrival for processing or does not have a new task arrival. This assumption is reasonable by setting the duration of each time slot to be small, e.g., $\Delta = 0.1$ second. Each mobile device has a scheduler. The scheduler places the newly arrived task to either a computation queue or a transmission queue (see Fig. 1) at the beginning of the next time slot. If the task is placed in the computation queue, then it will be processed locally. If the task is placed in the transmission queue, then it will be sent to an edge node through a wireless link for processing. Note that for the computation (or the transmission) queue, we assume that if the processing (or transmission) of a task is completed in a time slot, then the next task in the queue will be processed (or transmitted) at the beginning of the next time slot. This assumption is consistent with some existing works considering queuing dynamics in an MEC system (e.g., [21]), and the incurred additional delay can be ignored if the number of time slots that a task needs to wait in the queue and to be processed (or sent) is relatively large.

In the following, we first present the task model and the task offloading decision, respectively. Then, we introduce the computation and transmission queues.

2.1.1 Task Model

At the beginning of time slot $t \in \mathcal{T}$, if mobile device $m \in \mathcal{M}$ has a newly arrived task to be placed to a queue, then we define a variable $k_m(t) \in \mathbb{Z}_{++}$ to denote the unique index of the task. If mobile device m does not have a new task arrival to be placed at the beginning of time slot t , then $k_m(t)$ is set to zero for presentation simplicity.

Let $\lambda_m(t)$ (in bits) denote the number of newly arrived bits to be placed in a queue at the beginning of time slot $t \in \mathcal{T}$. If there exists a new task $k_m(t)$ at the beginning of time slot t , then $\lambda_m(t)$ is equal to the size of task $k_m(t)$. Otherwise, $\lambda_m(t)$ is set to zero. We set the size of task $k_m(t)$ to be from a discrete set $\Lambda \triangleq \{\lambda_1, \lambda_2, \dots, \lambda_{|\Lambda|}\}$ with $|\Lambda|$ available values. Hence, $\lambda_m(t) \in \Lambda \cup \{0\}$. In addition, task $k_m(t)$ requires a processing density of ρ_m (in CPU cycles per bit), i.e., the number of CPU cycles required to process a unit of data. Task $k_m(t)$ has a deadline τ_m (in time slots). That is, if task $k_m(t)$ has not been completely processed by the end of time slot $t + \tau_m - 1$, then it will be dropped immediately.

2.1.2 Task Offloading Decision

If mobile device $m \in \mathcal{M}$ has a newly arrived task $k_m(t)$ at the beginning of time slot $t \in \mathcal{T}$, then it needs to make two decisions for task $k_m(t)$. First, the mobile device decides whether to place the task to the computation queue or transmission queue. Second, if the task is placed to the transmission queue, then the mobile device decides the edge node to which it should offload the task.

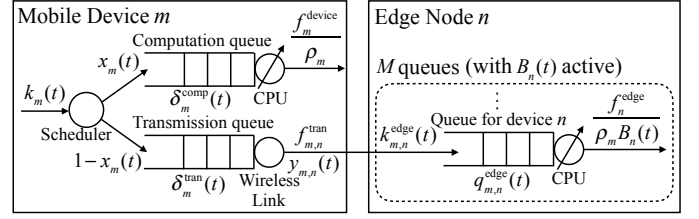


Fig. 1: An illustration of an MEC system with a mobile device $m \in \mathcal{M}$ and an edge node $n \in \mathcal{N}$.

Let binary variable $x_m(t) \in \{0, 1\}$ denote whether task $k_m(t)$ is scheduled to the computation queue or transmission queue. We set $x_m(t) = 1$ (or 0) if the task is scheduled to the computation queue (or the transmission queue). At the beginning of time slot t , $\lambda_m(t)x_m(t)$ is the number of bits arrived at the computation queue of mobile device m , and $\lambda_m(t)(1 - x_m(t))$ is the number of bits arrived at the transmission queue of mobile device m .

Let binary variable $y_{m,n}(t) \in \{0, 1\}$ denote whether task $k_m(t)$ is offloaded to edge node $n \in \mathcal{N}$ or not. We set $y_{m,n}(t) = 1$ if task $k_m(t)$ is offloaded to edge node n , and $y_{m,n}(t) = 0$ otherwise. For presentation convenience, we introduce vector $\mathbf{y}_m(t) = (y_{m,n}(t), n \in \mathcal{N})$. We assume that each task can be offloaded to only one edge node, i.e.,

$$\sum_{n \in \mathcal{N}} y_{m,n}(t) = \mathbb{1}(x_m(t) = 0), \quad m \in \mathcal{M}, t \in \mathcal{T}, \quad (1)$$

where the indicator $\mathbb{1}(z \in \mathcal{Z}) = 1$ if $z \in \mathcal{Z}$, and is equal to zero otherwise.

2.1.3 Computation Queue

For mobile device $m \in \mathcal{M}$, let f_m^{device} (in CPU cycles) denote its total processing capacity within each time slot. As a result, device m can process a maximum of $f_m^{\text{device}}/\rho_m$ bits for the computation queue within each time slot, i.e., the total processing capacity in each time slot divided by the required processing density of the tasks.

At the beginning of time slot $t \in \mathcal{T}$, if task $k_m(t)$ is placed in the computation queue, then we define a variable $l_m^{\text{comp}}(t) \in \mathcal{T}$ to denote the time slot when task $k_m(t)$ has either been processed or dropped. Without loss of generality, if either task $k_m(t)$ is not placed in the computation queue or $k_m(t) = 0$, then we set $l_m^{\text{comp}}(t) = 0$.

Let $\delta_m^{\text{comp}}(t)$ (in time slots) denote the number of time slots that task $k_m(t)$ will wait for processing if it is placed in the computation queue. Note that mobile device m will compute the value of $\delta_m^{\text{comp}}(t)$ before it decides the queue to place the task. Given $l_m^{\text{comp}}(t')$ for $t' < t$, the value of $\delta_m^{\text{comp}}(t)$ is computed as follows. For $m \in \mathcal{M}$ and $t \in \mathcal{T}$,

$$\delta_m^{\text{comp}}(t) = \left[\max_{t' \in \{0, 1, \dots, t-1\}} l_m^{\text{comp}}(t') - t + 1 \right]^+, \quad (2)$$

where the operator $[z]^+ = \max\{0, z\}$, and we set $l_m^{\text{comp}}(0) = 0$ for presentation simplicity. Specifically, the term $\max_{t' \in \{0, 1, 2, \dots, t-1\}} l_m^{\text{comp}}(t')$ determines the time slot when all the tasks placed in the computation queue before time slot t has either been processed or dropped. Hence, $\delta_m^{\text{comp}}(t)$ determines the number of time slots that task $k_m(t)$ should wait for processing. For example, suppose task $k_m(1)$ is placed in the computation queue, and its

processing will be completed in time slot 5, i.e., $l_m^{\text{comp}}(1) = 5$. In the meanwhile, suppose $k_m(2) = 0$, i.e., $l_m^{\text{comp}}(2) = 0$. At the beginning of time slot 3, if task $k_m(3)$ is placed in the computation queue, then its processing will start after time slot $l_m^{\text{comp}}(1) = 5$. Hence, it should wait for $\delta_m^{\text{comp}}(3) = [\max\{5, 0\} - 3 + 1]^+ = 3$ time slots.

If mobile device $m \in \mathcal{M}$ places task $k_m(t)$ in the computation queue at the beginning of time slot $t \in \mathcal{T}$ (i.e., $x_m(t) = 1$), then task $k_m(t)$ will have either been processed or dropped in time slot $l_m^{\text{comp}}(t)$:

$$l_m^{\text{comp}}(t) = \min \left\{ t + \delta_m^{\text{comp}}(t) + \left\lceil \frac{\lambda_m(t)}{f_m^{\text{device}}/\rho_m} \right\rceil - 1, t + \tau_m - 1 \right\}, \quad (3)$$

where $\lceil \cdot \rceil$ is the ceiling function. Specifically, the processing of task $k_m(t)$ will start at the beginning of time slot $t + \delta_m^{\text{comp}}(t)$. The number of time slots required to process the task is $\lceil \lambda_m(t)/(f_m^{\text{device}}/\rho_m) \rceil$. Hence, the first term in the min operator is the time slot when the processing of task $k_m(t)$ will be completed without considering the deadline of the task. The second term is the time slot when task $k_m(t)$ will be dropped. As a result, $l_m^{\text{comp}}(t)$ determines the time slot when task $k_m(t)$ will either be processed or dropped.

2.1.4 Transmission Queue

For mobile device $m \in \mathcal{M}$, let $f_{m,n}^{\text{tran}}$ (in bits) denote the total transmission capacity from the mobile device to edge node $n \in \mathcal{N}$ in each time slot.¹ At the beginning of time slot $t \in \mathcal{T}$, if task $k_m(t)$ is placed in the transmission queue, then we define a variable $l_m^{\text{tran}}(t) \in \mathcal{T}$ to denote the time slot when task $k_m(t)$ has been either sent or dropped. Without loss of generality, if either task $k_m(t)$ is not placed in the transmission queue or $k_m(t) = 0$, then we set $l_m^{\text{tran}}(t) = 0$.

Let $\delta_m^{\text{tran}}(t)$ (in time slots) denote the number of time slots that task $k_m(t)$ should wait for transmission if it is placed in the transmission queue. Note that mobile device m will compute the value of $\delta_m^{\text{tran}}(t)$ before it has decided on which queue to place the task. Given $l_m^{\text{tran}}(t')$ for $t' < t$, the value of $\delta_m^{\text{tran}}(t)$ is computed as follows. For $m \in \mathcal{M}$ and $t \in \mathcal{T}$,

$$\delta_m^{\text{tran}}(t) = \left[\max_{t' \in \{0, 1, \dots, t-1\}} l_m^{\text{tran}}(t') - t + 1 \right]^+, \quad (4)$$

where we set $l_m^{\text{tran}}(0) = 0$ for presentation simplicity.

If mobile device $m \in \mathcal{M}$ places task $k_m(t)$ in the transmission queue at the beginning of time slot $t \in \mathcal{T}$ (i.e., $x_m(t) = 0$), then task $k_m(t)$ will either be sent or dropped in time slot $l_m^{\text{tran}}(t)$:

$$l_m^{\text{tran}}(t) = \min \left\{ t + \delta_m^{\text{tran}}(t) + \left\lceil \sum_{n \in \mathcal{N}} \frac{y_{m,n}(t)\lambda_m(t)}{f_{m,n}^{\text{tran}}} \right\rceil - 1, t + \tau_m - 1 \right\}. \quad (5)$$

The idea of computing $l_m^{\text{tran}}(t)$ in (5) is similar as that of computing $l_m^{\text{comp}}(t)$ in (3).

1. Since we focus on characterizing the load level dynamics at the edge nodes, we consider a constant transmission capacity in the system model, as in some of the existing works such as [12], [14], [22], [23].

2.2 Edge Node Model

Each edge node $n \in \mathcal{N}$ maintains M queues, each queue corresponding to a mobile device in set \mathcal{M} . We assume that after an offloaded task is received by an edge node in a time slot, the task will be placed in its corresponding queue at the edge node at the beginning of the next time slot. This assumption is used to ensure that a task is processed by an edge node after the task has been completely received, and the incurred additional delay can be ignored if the number of time slots that each task needs for waiting and being processed is relatively large.

If a task of mobile device $m \in \mathcal{M}$ is placed in its corresponding queue at edge node $n \in \mathcal{N}$ at the beginning of time slot $t \in \mathcal{T}$, then we define a variable $k_{m,n}^{\text{edge}}(t) \in \mathbb{Z}_{++}$ to denote the unique index of the task.² Specifically, if task $k_m(t')$ for $t' \in \{1, 2, \dots, t-1\}$ is sent to edge node n in time slot $t-1$, then we have $k_{m,n}^{\text{edge}}(t) = k_m(t')$. Note that if there does not exist a task of mobile device m being placed in the queue at edge node n at the beginning of time slot t , then we set $k_{m,n}^{\text{edge}}(t) = 0$. Let $\lambda_{m,n}^{\text{edge}}(t) \in \Lambda \cup \{0\}$ (in bits) denote the number of bits arrived in the queue of mobile device m at edge node n at the beginning of time slot t . If task $k_{m,n}^{\text{edge}}(t)$ is placed in the corresponding queue at the beginning of time slot t , then $\lambda_{m,n}^{\text{edge}}(t)$ is equal to the size of task $k_{m,n}^{\text{edge}}(t)$. Otherwise, $\lambda_{m,n}^{\text{edge}}(t) = 0$.

Due to the unknown future load level dynamics at the edge nodes, mobile devices and edge nodes are unaware of the waiting and processing time of the tasks offloaded to the edge nodes until those tasks have either been processed or dropped. In the following, we first introduce the operation of the queues at the edge nodes. Then, we compute the time slot when a task has either been processed or dropped.

2.2.1 Queues at Edge Nodes

Let $q_{m,n}^{\text{edge}}(t)$ (in bits) denote the queue length of mobile device $m \in \mathcal{M}$ at edge node $n \in \mathcal{N}$ at the end of time slot $t \in \mathcal{T}$. Among those queues at edge node n , we refer to the queue of mobile device m as an *active queue* in time slot t if either there is a task of mobile device m arrived at the queue in time slot t (i.e., $\lambda_{m,n}^{\text{edge}}(t) > 0$) or the queue is non-empty at the end of time slot $t-1$ (i.e., $q_{m,n}^{\text{edge}}(t-1) > 0$). Let $\mathcal{B}_n(t)$ denote the set of active queues at edge node n in time slot t . That is, for $n \in \mathcal{N}$ and $t \in \mathcal{T}$,

$$\mathcal{B}_n(t) = \{m \mid m \in \mathcal{M}, \lambda_{m,n}^{\text{edge}}(t) > 0 \text{ or } q_{m,n}^{\text{edge}}(t-1) > 0\}. \quad (6)$$

Let $B_n(t)$ denote the number of active queues at edge node n in time slot t , i.e., $B_n(t) = |\mathcal{B}_n(t)|$.

Within time slot $t \in \mathcal{T}$, the active queues at an edge node $n \in \mathcal{N}$, i.e., the queues in set $\mathcal{B}_n(t)$, equally share the processing capacity of edge node n . This is a generalized processor sharing (GPS) model [24] with equal processing capacity sharing, and it can be approximated by practical algorithms such as the fair queuing algorithm in [25]. Let f_n^{edge} denote the total processing capacity of edge node n within each time slot. Therefore, edge node n can process

2. In each time slot, an edge node receives at most one task from a mobile device, as we have assumed that after a task is sent, the transmission of the next task starts at the beginning of the next time slot.

a maximum of $f_n^{\text{edge}}/(\rho_m B_n(t))$ bits for any active device $m \in \mathcal{B}_n(t)$ within time slot t .

For any queue at the edge nodes, we assume that if the processing of a task is completed in a time slot, then the next task in the queue will be processed at the beginning of the next time slot. To compute the queue length of mobile device $m \in \mathcal{M}$ at edge node $n \in \mathcal{N}$, let $e_{m,n}^{\text{edge}}(t)$ (in bits) denote the number of bits of the tasks dropped by the queue at the end of time slot $t \in \mathcal{T}$. Consequently, the queue length $q_{m,n}^{\text{edge}}(t)$ is updated as follows. For $m \in \mathcal{M}$, $n \in \mathcal{N}$, and $t \in \mathcal{T}$,

$$q_{m,n}^{\text{edge}}(t) = \left[q_{m,n}^{\text{edge}}(t-1) + \lambda_{m,n}^{\text{edge}}(t) - \frac{f_n^{\text{edge}}}{\rho_m B_n(t)} \mathbb{1}(m \in \mathcal{B}_n(t)) - e_{m,n}^{\text{edge}}(t) \right]^+ \quad (7)$$

Intuitively, the queue length $q_{m,n}^{\text{edge}}(t)$ is equal to the queue length in the previous time slot $q_{m,n}^{\text{edge}}(t-1)$ plus the difference between the bits arrived in time slot t and the bits being served (i.e., processed or dropped) in time slot t .

2.2.2 Task Processing or Dropping

If task $k_{m,n}^{\text{edge}}(t)$ of mobile device $m \in \mathcal{M}$ is placed in the corresponding queue at edge node $n \in \mathcal{N}$ at the beginning of time slot $t \in \mathcal{T}$, then we define a variable $l_{m,n}^{\text{edge}}(t) \in \mathcal{T}$ to denote the time slot when this task has either been processed or dropped by edge node n . Due to the uncertain future load at edge node n , the value of $l_{m,n}^{\text{edge}}(t)$ is unknown to mobile device m and edge node n until the associated task $k_{m,n}^{\text{edge}}(t)$ has either been processed or dropped. Without loss of generality, if $k_{m,n}^{\text{edge}}(t) = 0$, then we set $l_{m,n}^{\text{edge}}(t) = 0$.

For the definition of variable $l_{m,n}^{\text{edge}}(t)$, let $\hat{l}_{m,n}^{\text{edge}}(t)$ denote the time slot when the processing of task $k_{m,n}^{\text{edge}}(t)$ starts, i.e., for $m \in \mathcal{M}$, $n \in \mathcal{N}$, and $t \in \mathcal{T}$,

$$\hat{l}_{m,n}^{\text{edge}}(t) = \max \left\{ t, \max_{t' \in \{0,1,\dots,t-1\}} l_{m,n}^{\text{edge}}(t') + 1 \right\}, \quad (8)$$

where we set $l_{m,n}^{\text{edge}}(0) = 0$. Specifically, the time slot when the processing of task $k_{m,n}^{\text{edge}}(t)$ starts should be no earlier than the time slot that the task is placed in the queue or the time slot when each of the tasks arrived earlier has been processed or dropped.

Given the realization of the load levels at edge node n , $l_{m,n}^{\text{edge}}(t)$ is the time slot satisfying the following constraints. For $m \in \mathcal{M}$, $n \in \mathcal{N}$, and $t \in \mathcal{T}$,

$$\sum_{t'=\hat{l}_{m,n}^{\text{edge}}(t)}^{l_{m,n}^{\text{edge}}(t)} \frac{f_n^{\text{edge}}}{\rho_m B_n(t')} \mathbb{1}(m \in \mathcal{B}_n(t')) \geq \lambda_{m,n}^{\text{edge}}(t), \quad (9)$$

$$\sum_{t'=\hat{l}_{m,n}^{\text{edge}}(t)}^{l_{m,n}^{\text{edge}}(t)-1} \frac{f_n^{\text{edge}}}{\rho_m B_n(t')} \mathbb{1}(m \in \mathcal{B}_n(t')) < \lambda_{m,n}^{\text{edge}}(t). \quad (10)$$

Specifically, the total processing capacity that edge node n allocated to mobile device m from time slot $\hat{l}_{m,n}^{\text{edge}}(t)$ to time slot $l_{m,n}^{\text{edge}}(t)$ should be no smaller than the size of task $k_{m,n}^{\text{edge}}(t)$, while the corresponding total processing capacity allocated from time slot $\hat{l}_{m,n}^{\text{edge}}(t)$ to time slot $l_{m,n}^{\text{edge}}(t) - 1$ should be smaller than the size of the task.

3 TASK OFFLOADING PROBLEM IN MEC

In this section, we present the task offloading problem for the MEC system. Specifically, at the beginning of each time slot, each mobile device observes its state (e.g., task size, queue information). If there is a newly arrived task to be processed, then the mobile device chooses an action for the task. The observed state and the chosen action will result in a cost (i.e., the delay of the task if the task is processed, or a penalty if it is dropped) for the mobile device. The objective of each mobile device is to minimize its expected long-term cost by optimizing the policy mapping from states to actions. In the following, we first introduce the state, action, and cost, respectively. We then formulate the cost minimization problem for each device.

3.1 State

At the beginning of time slot $t \in \mathcal{T}$, each device $m \in \mathcal{M}$ observes its state information, including the task size, the information related to the queues, and the load level history at the edge nodes. Specifically, mobile device m maintains the following state vector:

$$s_m(t) = \left(\lambda_m(t), \delta_m^{\text{comp}}(t), \delta_m^{\text{tran}}(t), \mathbf{q}_m^{\text{edge}}(t-1), \mathbf{H}(t) \right), \quad (11)$$

where vector $\mathbf{q}_m^{\text{edge}}(t-1) = (q_{m,n}^{\text{edge}}(t-1), n \in \mathcal{N})$. The matrix $\mathbf{H}(t)$ includes the history of the load level (i.e., the number of active queues) of each edge node within the previous T^{step} time slots (i.e., from time slot $t - T^{\text{step}}$ to time slot $t - 1$), with which the load levels at the edge nodes in the near future can be estimated. It is a matrix with size $T^{\text{step}} \times N$. Let $\{\mathbf{H}(t)\}_{i,j}$ denote the (i,j) element of matrix $\mathbf{H}(t)$, and it corresponds to the load level history of edge node j in the i^{th} time slot starting from time slot $t - T^{\text{step}}$, i.e., time slot $t - T^{\text{step}} + i - 1$. The element $\{\mathbf{H}(t)\}_{i,j}$ is defined as follows:

$$\{\mathbf{H}(t)\}_{i,j} = B_j(t - T^{\text{step}} + i - 1), \quad (12)$$

which is the number of active queues of edge node j in time slot $t - T^{\text{step}} + i - 1$. Let \mathcal{S} denote the discrete and finite state space of each mobile device. Formally, set $\mathcal{S} = \Lambda \times \{0, 1, \dots, T\}^2 \times \mathcal{Q}^N \times \{0, 1, \dots, M\}^{T^{\text{step}} \times N}$, where \mathcal{Q} denotes the set of the available values of the queue length at an edge node within the T time slots.

Mobile device $m \in \mathcal{M}$ can obtain state information $\lambda_m(t)$, $\delta_m^{\text{comp}}(t)$, and $\delta_m^{\text{tran}}(t)$ through local observation at the beginning of time slot t . For state information $\mathbf{q}_m^{\text{edge}}(t-1)$, mobile device m can compute this vector according to the number of bits of device m transmitted to each edge node in each time slot and the number of bits of device m processed or being dropped by each edge node in each time slot according to (7). For matrix $\mathbf{H}(t)$, we assume that each edge node will broadcast its number of active queues at the end of each time slot. Since the number of active queues is always a small number, which can be represented by several bits, the broadcasting will only incur a small signaling overhead.

3.2 Action

At the beginning of time slot $t \in \mathcal{T}$, if mobile device $m \in \mathcal{N}$ has a new task arrival $k_m(t)$, then it will choose the actions for task $k_m(t)$: (a) whether to schedule the task to the computation queue or the transmission queue, i.e.,

$x_m(t)$; (b) which edge node the task is offloaded to, i.e., $\mathbf{y}_m(t) = (y_{m,n}(t), n \in \mathcal{N})$. Hence, the action of device m in time slot t is represented by the following action vector:

$$\mathbf{a}_m(t) = (x_m(t), \mathbf{y}_m(t)). \quad (13)$$

Let \mathcal{A} denote the decision space of each mobile device, i.e., $\mathcal{A} = \{0, 1\}^{1+N}$.

3.3 Cost

Let $d_m(\mathbf{s}_m(t), \mathbf{a}_m(t))$ (in time slots) denote the delay of task $k_m(t)$, given the observed state $\mathbf{s}_m(t)$ and the selected action $\mathbf{a}_m(t)$. For $m \in \mathcal{M}$ and $t \in \mathcal{T}$, if $x_m(t) = 1$, then

$$d_m(\mathbf{s}_m(t), \mathbf{a}_m(t)) = l_m^{\text{comp}}(t) - t + 1; \quad (14)$$

if $x_m(t) = 0$, then

$$\begin{aligned} d_m(\mathbf{s}_m(t), \mathbf{a}_m(t)) &= \sum_{n \in \mathcal{N}} \sum_{t'=t}^T \mathbb{1}(k_{m,n}^{\text{edge}}(t') = k_m(t)) l_{m,n}^{\text{edge}}(t') - t + 1. \end{aligned} \quad (15)$$

Specifically, the delay of task $k_m(t)$ is the number of time slots between time slot t and the time slot when task $k_m(t)$ has either been processed or dropped.

There is a cost $c_m(\mathbf{s}_m(t), \mathbf{a}_m(t))$ associated with task $k_m(t)$. If task $k_m(t)$ has been processed, then

$$c_m(\mathbf{s}_m(t), \mathbf{a}_m(t)) = d_m(\mathbf{s}_m(t), \mathbf{a}_m(t)). \quad (16)$$

On the other hand, if task $k_m(t)$ has been dropped, then

$$c_m(\mathbf{s}_m(t), \mathbf{a}_m(t)) = C, \quad (17)$$

where $C > 0$ is a constant penalty. Without loss of generality, if task $k_m(t) = 0$, then we set $c_m(\mathbf{s}_m(t), \mathbf{a}_m(t)) = 0$. In the remaining part of this work, we use the short form $c_m(t)$ to denote $c_m(\mathbf{s}_m(t), \mathbf{a}_m(t))$.

3.4 Problem Formulation

A policy of device $m \in \mathcal{M}$ is a mapping from its state to its action, i.e., $\pi_m : \mathcal{S} \rightarrow \mathcal{A}$. We aim to find the optimal policy π_m^* for each device m such that its expected long-term cost is minimized, i.e.,

$$\begin{aligned} \pi_m^* &= \arg \underset{\pi_m}{\text{minimize}} \quad \mathbb{E} \left[\sum_{t \in \mathcal{T}} \gamma^{t-1} c_m(t) \mid \pi_m \right] \\ &\text{subject to} \quad \text{constraints (1) – (5), (7) – (10),} \\ &\quad \quad \quad (14) – (17), \end{aligned} \quad (18)$$

where $\gamma \in (0, 1]$ is a discount factor that characterizes the discounted cost in the future. The expectation $\mathbb{E}[\cdot]$ is with respect to the time-varying system environments, including the task arrivals and the computational requirements of the tasks of all mobile devices as well as the offloading decisions of the mobile devices other than device m .

Solving problem (18) is challenging. This is mainly due to the unknown load levels at the edge nodes, which depend on the decisions and the task models (e.g., the size and arrival time of the task) of other mobile devices, as well as the unknown future task models of the device itself. In this work, we propose a DRL-based offloading algorithm that addresses the challenge by learning the mapping from each state-action pair to its expected long-term cost.

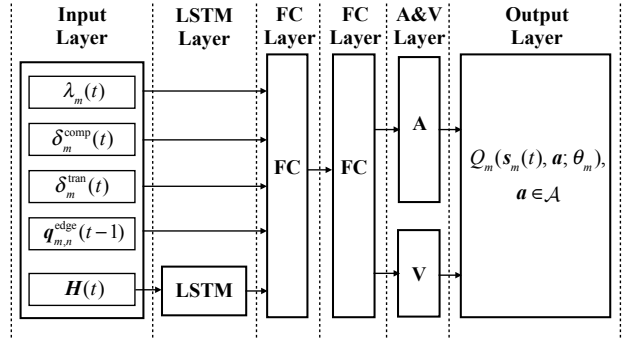


Fig. 2: The neural network of mobile device $m \in \mathcal{M}$ with parameter vector θ_m , which maps from state $\mathbf{s}_m(t) \in \mathcal{S}$ to the Q-value of each action $\mathbf{a} \in \mathcal{A}$.

4 DRL-BASED OFFLOADING ALGORITHM

In this section, we propose a DRL-based offloading algorithm that enables the distributed offloading decision making of each mobile device. In the proposed algorithm, each mobile device aims to learn a mapping from each state-action pair to a Q-value, which characterizes the expected long-term cost of the state-action pair. The mapping is determined by a neural network. Based on the mapping, each device can select the action inducing the minimum Q-value under its state to minimize its expected long-term cost.

In the following, we first introduce the neural network for a mobile device that characterizes its mapping from state-action pairs to Q-values. Then, we present the DRL-based algorithm and describe the message exchange between a mobile device and an edge node.

4.1 Neural Network

The objective of the neural network is to find a mapping from each state to a set of Q-values, each corresponding to an action. As shown in Fig. 2, for any mobile device $m \in \mathcal{M}$, we consider a neural network with six layers: an input layer, an LSTM layer, two fully connected (FC) layers, an advantage and value (A&V) layer, and an output layer. Let θ_m denote the parameter vector of the neural network of device m , which includes the weights of all connections and the biases of all neurons from the input layer to the A&V layer.³ The details of each layer are as follows.

4.1.1 Input Layer

This layer is responsible for taking the state vector as input and passing them to the following layers. For mobile device $m \in \mathcal{M}$, the state information includes $\lambda_m(t)$, $\delta_m^{\text{comp}}(t)$, $\delta_m^{\text{tran}}(t)$, $\mathbf{q}_m^{\text{edge}}(t-1)$, and $\mathbf{H}(t)$. The state information $\lambda_m(t)$, $\delta_m^{\text{comp}}(t)$, $\delta_m^{\text{tran}}(t)$, and $\mathbf{q}_m^{\text{edge}}(t-1)$ will be passed to the FC layer, and $\mathbf{H}(t)$ will be passed to the LSTM layer.

4.1.2 LSTM Layer

This layer is responsible for learning the dynamics of the load levels at the edge nodes. This is achieved by including

3. The weights of the connections between the A&V layer and the output layer as well as the bias of the neurons in the output layer are given and fixed. Hence, we do not include them in the network parameter vector θ_m , as the vector θ_m includes the parameters that are adjustable through learning in the DRL-based algorithm.

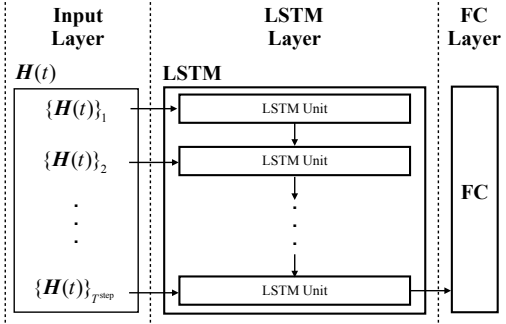


Fig. 3: An LSTM network with T^{step} LSTM units.

an LSTM network [26], [27]. We use the LSTM network because it can keep track of the state $\mathbf{H}(t)$ over time. It can provide the neural network the ability of estimating the load levels at the edge nodes in the future using the history.

Specifically, the LSTM network takes the matrix $\mathbf{H}(t)$ as input so as to learn the load level dynamics. Fig. 3 shows the structure of an LSTM network. The LSTM network contains T^{step} LSTM units, each of which contains a set of hidden neurons. Each LSTM unit takes one row of the matrix $\mathbf{H}(t)$ as input, we let $\{\mathbf{H}(t)\}_i$ denote the i^{th} row of matrix $\mathbf{H}(t)$ in Fig. 3. These LSTM units are connected in sequence so as to keep track of the variations of the sequences from $\{\mathbf{H}(t)\}_1$ to $\{\mathbf{H}(t)\}_{T^{\text{step}}}$, which can reveal the variations of the load levels at the edge nodes among time slots. The LSTM network will output the information that indicates the dynamics of the load levels in the future in the last LSTM unit, where the output will be connected to the neurons in the next layer for further learning.

4.1.3 FC Layers

The two FC layers are responsible for learning the mapping from the state and the learned load level dynamics to the Q-values of the actions. Each FC layer contains a set of neurons with rectified linear unit (ReLU). For the first FC layer, the input of each neuron connects to the neurons in the input layer corresponding to all the states (except the matrix $\mathbf{H}(t)$) and the LSTM network in the LSTM layer. The output of each neuron connects to each of the neurons in the second FC layer. For the second FC layer, the output of each neuron connects to each of the neurons in the A&V layer.

4.1.4 A&V Layer and Output Layer

The A&V layer and the output layer implement the dueling-DQN technique [28] and determine the Q-value of each action as output. The main idea of the dueling-DQN is to first separately learn a *state-value* (i.e., the portion of the Q-value resulting from the state) and *action-advantage values* (i.e., the portion of the Q-value resulting from the actions), and then use the state-value and action-advantage values to determine the Q-values of state-action pairs. This technique can improve the estimation of the Q-values through separately evaluating the expected long-term cost resulting from a state and an action.

The A&V layer contains two networks, denoted by network A and network V (see Fig. 2). The network A contains an FC network with a set of neurons. It is responsible for learning the action-advantage value of each action $\mathbf{a} \in \mathcal{A}$.

For mobile device $m \in \mathcal{M}$, let $A_m(\mathbf{s}_m(t), \mathbf{a}; \boldsymbol{\theta}_m)$ denote the action-advantage value of action \mathbf{a} under state $\mathbf{s}_m(t) \in \mathcal{S}$ with network parameter vector $\boldsymbol{\theta}_m$. The network V contains an FC network with a set of neurons. It is responsible for learning the state-value. For mobile device m , let $V_m(\mathbf{s}_m(t); \boldsymbol{\theta}_m)$ denote state-value of state $\mathbf{s}_m(t)$ with network parameter vector $\boldsymbol{\theta}_m$. The values of $A_m(\mathbf{s}_m(t), \mathbf{a}; \boldsymbol{\theta}_m)$ and $V_m(\mathbf{s}_m(t); \boldsymbol{\theta}_m)$ are determined by the parameter vector $\boldsymbol{\theta}_m$ and the neural network structure from the input layer to the A&V layer, where vector $\boldsymbol{\theta}_m$ is adjustable and will be trained in the DRL-based algorithm.

Based on the A&V layer, for mobile device $m \in \mathcal{M}$, the resulting Q-value of action $\mathbf{a} \in \mathcal{A}$ under state $\mathbf{s}_m(t) \in \mathcal{S}$ in the output layer is given as follows [28]:

$$Q_m(\mathbf{s}_m(t), \mathbf{a}; \boldsymbol{\theta}_m) = V_m(\mathbf{s}_m(t); \boldsymbol{\theta}_m) + \left(A_m(\mathbf{s}_m(t), \mathbf{a}; \boldsymbol{\theta}_m) - \frac{1}{|\mathcal{A}|} \sum_{\mathbf{a}' \in \mathcal{A}} A_m(\mathbf{s}_m(t), \mathbf{a}'; \boldsymbol{\theta}_m) \right), \quad (19)$$

which is the sum of the state-value under the corresponding state and the additional action-advantage value of the corresponding action (i.e., the difference between the action-advantage value of the action and the average action-advantage value over all actions).

In summary, from the input layer to the output layer, the neural network of mobile device $m \in \mathcal{M}$ with parameter vector $\boldsymbol{\theta}_m$ forms a mapping from state-action pairs to Q-values (i.e., under any observed state $\mathbf{s}_m(t) \in \mathcal{S}$, there is a Q-value for each action $\mathbf{a} \in \mathcal{A}$, denoted by $Q_m(\mathbf{s}_m(t), \mathbf{a}; \boldsymbol{\theta}_m)$), which characterizes the expected long-term cost under the observed state and each of the actions in the action space.

4.2 DRL-Based Algorithm

In our proposed DRL-based algorithm, we let edge nodes help mobile devices to train the neural network to alleviate the computational loads at the mobile devices. Specifically, for each mobile device $m \in \mathcal{M}$, there is an edge node $n_m \in \mathcal{N}$ which helps device m with the training. This edge node n_m can be the edge node that has the maximum transmission capacity with mobile device m . For presentation convenience, let $\mathcal{M}_n \subset \mathcal{M}$ denote the set of mobile devices whose training is performed by edge node $n \in \mathcal{N}$, i.e., $\mathcal{M}_n = \{m \in \mathcal{M} \mid n_m = n\}$. Note that it is reasonable to let edge nodes help with the training directly. This is because the information exchange involved in the training (including the state information and the neural network parameter) is small. In addition, the required processing capacity of the training in each time slot can be much less than those of the tasks of mobile devices.

The DRL-based algorithm to be executed at mobile device $m \in \mathcal{M}$ and edge node $n \in \mathcal{N}$ are given in Algorithms 1 and 2, respectively. The key idea of the algorithm is to train the neural network using the experience⁴ (i.e., state, action, cost, and next state) of the mobile device to obtain

4. We use the term “experience” to refer to the tuple consisting of state, action, cost, and next state, as used in [28], [29]. Alternatively, some other existing works, such as [16], [30], used the term “transition”.

Algorithm 1 DRL-based Algorithm at Device $m \in \mathcal{M}$

```

1: for episode from 1 to #_of_Episodes do
2:   Initialize  $\mathbf{s}_m(1)$ ;
3:   for time slot  $t \in \mathcal{T}$  do
4:     if device  $m$  has a new task arrival  $k_m(t)$  then
5:       Send a parameter_request to edge node  $n_m$ ;
6:       Receive network parameter vector  $\theta_m^{\text{Eval}}$ ;
7:       Select an action  $\mathbf{a}_m(t)$  according to (21);
8:     end if
9:     Observe the next state  $\mathbf{s}_m(t+1)$ ;
10:    Observe a set of costs  $\{c_m(t'), t' \in \tilde{\mathcal{T}}_{m,t}\}$ ;
11:    for each task  $k_m(t')$  with  $t' \in \tilde{\mathcal{T}}_{m,t}$  do
12:      Send  $(\mathbf{s}_m(t'), \mathbf{a}_m(t'), c_m(t'), \mathbf{s}_m(t'+1))$  to  $n_m$ ;
13:    end for
14:  end for
15: end for

```

the mapping from state-action pairs to Q-values, based on which the device can select the action leading to the minimum Q-value under the observed state to minimize its expected long-term cost.

In the DRL-based algorithm, the edge node $n \in \mathcal{N}$ maintains a replay memory D_m and two neural networks for device $m \in \mathcal{M}_n$. The replay memory D_m stores the observed experience $(\mathbf{s}_m(t), \mathbf{a}_m(t), c_m(t), \mathbf{s}_m(t+1))$ of mobile device m for some $t \in \mathcal{T}$, where we refer $(\mathbf{s}_m(t), \mathbf{a}_m(t), c_m(t), \mathbf{s}_m(t+1))$ as experience t of mobile device m . The experience in the replay memory is used to train the neural networks. The two neural networks include an *Eval_Net_m* and a *Target_Net_m*, and their Q-values are represented by $Q_m^{\text{Eval}}(\mathbf{s}_m(t), \mathbf{a}; \theta_m^{\text{Eval}})$ and $Q_m^{\text{Target}}(\mathbf{s}_m(t), \mathbf{a}; \theta_m^{\text{Target}})$ under observed state $\mathbf{s}_m(t) \in \mathcal{S}$ and action $\mathbf{a} \in \mathcal{A}$, respectively. Note that the *Eval_Net_m* and the *Target_Net_m* have the same neural network structure, as presented in Section 4.1, while they have different network parameter vectors θ_m^{Eval} and θ_m^{Target} , respectively. The *Eval_Net_m* is used for action selection. The *Target_Net_m* is used for characterizing a target Q-value, which approximates the expected long-term cost of an action under the observed state. This target Q-value will be used for updating the network parameter vector θ_m^{Eval} in *Eval_Net_m* by minimizing the difference between the Q-value under *Eval_Net_m* and the target Q-value. The initialization of the replay memory D_m and two neural networks are given in steps 1–3 in Algorithm 2.

In the following, we present the DRL-based algorithm at mobile device $m \in \mathcal{M}$ and edge node $n \in \mathcal{N}$, respectively.

4.2.1 Algorithm 1 at Mobile Device $m \in \mathcal{M}$

We consider multiple episodes, where #_of_Episodes denotes the number of episodes. At the beginning of each episode, mobile device $m \in \mathcal{M}$ initializes the state, i.e.,

$$\mathbf{s}_m(1) = (\lambda_m(1), \delta_m^{\text{comp}}(1), \delta_m^{\text{tran}}(1), \mathbf{q}_m^{\text{edge}}(0), \mathbf{H}(1)), \quad (20)$$

where we set $q_{m,n}^{\text{edge}}(0) = 0$ for all $n \in \mathcal{N}$, and $\mathbf{H}(1)$ is a zero matrix with size $T^{\text{step}} \times N$.⁵ Each episode contains a set of time slots \mathcal{T} .

At the beginning of time slot $t \in \mathcal{T}$, if mobile device m has a new task arrival $k_m(t)$, then it will send a *parameter_request* to edge node n_m . Upon receiving the requested

5. For matrix $\mathbf{H}(t)$, for any $t-1 < T^{\text{step}}$ (i.e., the number of observed history is smaller than T^{step}), $\{\mathbf{H}(t)\}_i = \mathbf{0}$ for $i = 1, 2, \dots, T^{\text{step}} - (t-1)$, where $\mathbf{0}$ is a zero vector with size N .

Algorithm 2 DRL-Based Algorithm at Edge Node $n \in \mathcal{N}$

```

1: Initialize replay memory  $D_m$  for each device  $m \in \mathcal{M}_n$  and
   set Count := 0;
2: Initialize Eval_Netm with random parameter  $\theta_m^{\text{Eval}}$  for each
   device  $m \in \mathcal{M}_n$ ;
3: Initialize Target_Netm with random parameter  $\theta_m^{\text{Target}}$  for
   each device  $m \in \mathcal{M}_n$ ;
4: while True do
5:   if receive a parameter_request from device  $m \in \mathcal{M}_n$  then
6:     Send  $\theta_m^{\text{Eval}}$  to device  $m$ ;
7:   end if
8:   if receive an experience  $(\mathbf{s}_m(t), \mathbf{a}_m(t), c_m(t), \mathbf{s}_m(t+1))$ 
   from device  $m \in \mathcal{M}_n$  then
9:     Store  $(\mathbf{s}_m(t), \mathbf{a}_m(t), c_m(t), \mathbf{s}_m(t+1))$  in  $D_m$ ;
10:    Sample a set of experiences (denoted by  $\mathcal{I}$ ) from  $D_m$ ;
11:    for each experience  $i \in \mathcal{I}$  do
12:      Obtain experience  $(\mathbf{s}_m(i), \mathbf{a}_m(i), c_m(i), \mathbf{s}_m(i+1))$ ;
13:      Compute  $\hat{Q}_{m,i}^{\text{Target}}$  according to (24);
14:    end for
15:    Set vector  $\hat{Q}_m^{\text{Target}} := (\hat{Q}_{m,i}^{\text{Target}}, i \in \mathcal{I})$ ;
16:    Update  $\theta_m^{\text{Eval}}$  to minimize  $L(\theta_m^{\text{Eval}}, \hat{Q}_m^{\text{Target}})$  in (23);
17:    Count := Count + 1;
18:    if mod(Count, Replace_Threshold) = 0 then
19:       $\theta_m^{\text{Target}} := \theta_m^{\text{Eval}}$ ;
20:    end if
21:  end if
22: end while

```

parameter vector θ_m^{Eval} of *Eval_Net_m* from edge node n_m , device m will choose its action for task $k_m(t)$ as follows:

$$\mathbf{a}_m(t) = \begin{cases} \text{select a random action from } \mathcal{A}, & \text{with prob. } \epsilon, \\ \arg \min_{\mathbf{a} \in \mathcal{A}} Q_m^{\text{Eval}}(\mathbf{s}_m(t), \mathbf{a}; \theta_m^{\text{Eval}}), & \text{with prob. } 1 - \epsilon, \end{cases} \quad (21)$$

where ‘prob.’ is the short-form for probability, and ϵ is the probability of random exploration. The value of $Q_m^{\text{Eval}}(\mathbf{s}_m(t), \mathbf{a}; \theta_m^{\text{Eval}})$ is the Q-value under the current parameter θ_m^{Eval} of neural network *Eval_Net_m*. Intuitively, with a probability $1 - \epsilon$, the mobile device chooses the action that corresponds to the minimum Q-value under the observed state $\mathbf{s}_m(t)$ based on *Eval_Net_m*.

At the beginning of the next time slot (i.e., time slot $t+1$), mobile device m observes the next state $\mathbf{s}_m(t+1)$. On the other hand, as the processing and the transmission of a task may continue for multiple time slots, the cost $c_m(t)$, which depends on the delay of task $k_m(t)$, may not be observed at the beginning of time slot $t+1$. Instead, mobile device m may observe a set of costs belonging to some tasks $k_m(t')$ with time slot $t' \leq t$. To address this, for device m , we define $\tilde{\mathcal{T}}_{m,t} \subset \mathcal{T}$ as the set of time slots such that each task $k_m(t')$ associated with time slot $t' \in \tilde{\mathcal{T}}_{m,t}$ has been processed or dropped in time slot t . Set $\tilde{\mathcal{T}}_{m,t}$ is defined as follows:

$$\tilde{\mathcal{T}}_{m,t} = \left\{ t' \mid t' = 1, 2, \dots, t, \lambda_m(t') > 0, x_m(t') l_m^{\text{comp}}(t') + (1 - x_m(t')) \sum_{n \in \mathcal{N}} \sum_{i=t'}^t \mathbb{1}(k_{m,n}^{\text{edge}}(i) = k_m(t')) l_{m,n}^{\text{edge}}(i) = t \right\}. \quad (22)$$

In (22), $\lambda_m(t') > 0$ implies that there is a newly arrived task $k_m(t')$ in time slot t' . Specifically, set $\tilde{\mathcal{T}}_{m,t}$ contains a time

slot $t' \in \{1, 2, \dots, t\}$ if task $k_m(t')$ has been processed or dropped in time slot t . Hence, at the beginning of time slot $t+1$, mobile device m can observe a set of costs $\{c_m(t'), t' \in \tilde{\mathcal{T}}_{m,t}\}$, where set $\tilde{\mathcal{T}}_{m,t}$ can be an empty set for some $m \in \mathcal{M}$ and $t \in \mathcal{T}$. Then, for each task $k_m(t')$ with $t' \in \tilde{\mathcal{T}}_{m,t}$, device m sends its experience $(\mathbf{s}_m(t'), \mathbf{a}_m(t'), c_m(t'), \mathbf{s}_m(t'+1))$ to edge node n_m .

4.2.2 Algorithm 2 at Edge Node $n \in \mathcal{N}$

After initializing the replay memory D_m as well as the neural networks $Eval_Net_m$ and $Target_Net_m$ for device $m \in \mathcal{M}_n$, edge node $n \in \mathcal{N}$ will wait for the request messages from the mobile devices in set \mathcal{M}_n . If edge node n receives a *parameter_request* from mobile device $m \in \mathcal{M}_n$, then it will send the current parameter vector θ_m^{Eval} of $Eval_Net_m$ to device m . On the other hand, if edge node n receives an experience $(\mathbf{s}_m(t), \mathbf{a}_m(t), c_m(t), \mathbf{s}_m(t+1))$ from mobile device $m \in \mathcal{M}_n$, then it will store the experience in memory D_m . The memory has a maximum size, and it serves in a first-in first-out (FIFO) manner. Note that we do not require synchronization between mobile device m and its associated edge node n_m . The edge node will train the neural network (in steps 10–20 in Algorithm 2) to update the parameter vector θ_m^{Eval} of $Eval_Net_m$ as follows.

The edge node will randomly sample a set of experiences from the memory (in step 10), denoted by \mathcal{I} . Based on these experience samples, the key idea of the update of $Eval_Net_m$ is to minimize the difference between the Q-values under $Eval_Net_m$ and the target Q-values computed based on the experience samples under $Target_Net_m$. Specifically, for the experience samples in set \mathcal{I} , the edge node will compute $\hat{Q}_m^{\text{Target}} = (\hat{Q}_{m,i}^{\text{Target}}, i \in \mathcal{I})$ and update θ_m^{Eval} in $Eval_Net_m$ by minimizing the following loss function:

$$L(\theta_m^{\text{Eval}}, \hat{Q}_m^{\text{Target}}) = \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \left(Q_m^{\text{Eval}}(\mathbf{s}_m(i), \mathbf{a}_m(i); \theta_m^{\text{Eval}}) - \hat{Q}_{m,i}^{\text{Target}} \right)^2, \quad (23)$$

where $|\mathcal{I}|$ is the cardinality of set \mathcal{I} . The loss function (23) characterizes the gap between the Q-value of action $\mathbf{a}_m(i)$ given state $\mathbf{s}_m(i)$ under the current network parameter vector θ_m^{Eval} and a target Q-value $\hat{Q}_{m,i}^{\text{Target}}$ for each experience $i \in \mathcal{I}$ (to be explained in the next paragraph). The minimization of the loss function is accomplished by performing a gradient descent step on the neural network $Eval_Net_m$ using backpropagation (see Section 6 in [31]).

The value of $\hat{Q}_{m,i}^{\text{Target}}$ for experience $i \in \mathcal{I}$ is determined based on double-DQN technique [32]. The double-DQN technique can improve the estimation of the expected long-term cost when compared with the traditional method (e.g., [16]). The value of $\hat{Q}_{m,i}^{\text{Target}}$ for experience i is the sum of the corresponding cost in experience i and a discounted Q-value of the action that is likely to be selected given the next state in experience i under network $Target_Net_m$, i.e.,

$$\hat{Q}_{m,i}^{\text{Target}} = c_m(i) + \gamma Q_m^{\text{Target}}(\mathbf{s}_m(i+1), \mathbf{a}_i^{\text{Next}}; \theta_m^{\text{Target}}), \quad (24)$$

where $\mathbf{a}_i^{\text{Next}}$ is the action with the minimum Q-value given state $\mathbf{s}_m(i+1)$ under $Eval_Net_m$, i.e.,

$$\mathbf{a}_i^{\text{Next}} = \arg \min_{\mathbf{a} \in \mathcal{A}} Q_m^{\text{Eval}}(\mathbf{s}_m(i+1), \mathbf{a}; \theta_m^{\text{Eval}}). \quad (25)$$

TABLE 1: Parameter settings

Parameter	Value
M	50
N	5
Δ	0.1 second
$f_m^{\text{device}}, m \in \mathcal{M}$	2.5 GHz [15]
$f_n^{\text{edge}}, n \in \mathcal{N}$	41.8 GHz [15]
$f_{m,n}^{\text{tran}}, m \in \mathcal{M}, n \in \mathcal{N}$	14 Mbps [33]
$\lambda_m(t), m \in \mathcal{M}, t \in \mathcal{T}$	{2.0, 2.1, ..., 5.0} Mbits [4]
$\rho_m, m \in \mathcal{M}$	0.297 gigacycles per Mbits [4]
$\tau_m, m \in \mathcal{M}$	10 time slots (i.e., 1 second)
Task arrival probability	0.3

Intuitively, for experience i , the target-Q value $\hat{Q}_{m,i}^{\text{Target}}$ reveals the expected long-term cost of action $\mathbf{a}_m(i)$ given state $\mathbf{s}_m(i)$. This is the summation of the actual cost recorded in experience i , i.e., $c_m(i)$, and the approximate expected long-term future cost based on $Target_Net_m$, i.e., $\gamma Q_m^{\text{Target}}(\mathbf{s}_m(i+1), \mathbf{a}_i^{\text{Next}}; \theta_m^{\text{Target}})$.

For every Replace_Threshold updates, $Target_Net_m$ will be updated by copying $Eval_Net_m$, i.e., $\theta_m^{\text{Target}} = \theta_m^{\text{Eval}}$, where $\text{mod}(\cdot)$ is the modulo operator (in step 18 in Algorithm 2). The objective of this step is to keep the network parameter θ_m^{Target} in $Target_Net_m$ up-to-date, so that it can better approximate the expected long-term cost in the computing of the target Q-values in (24).

5 PERFORMANCE EVALUATION

In this section, we compare our proposed DRL-based method with several benchmark methods, including no offloading (denoted by No Offl.), random offloading (denoted by R. Offl.), PGOA in [14], and ULOOF in [15]. The PGOA is designed based on the best response algorithm for the potential game, which takes into account the strategic interaction among mobile devices. The ULOOF is designed based on the capacity estimation according to historical observations. In these simulations, we consider two performance metrics: the ratio of dropped tasks (i.e., the ratio of the number of dropped tasks to the number of total task arrivals) and the average delay (i.e., the average delay of the tasks whose processing has been completed). Unless stated otherwise, the basic parameter setting in the simulations are given in Table 1. In addition, the probability of random exploration ϵ is set to be gradually decreasing from 1 to 0.01, and the discount factor γ is set to be 0.9.

In the following, we first show the convergence of the proposed algorithm across episodes. Then, we compare the performance of our proposed algorithm with the benchmark methods under different parameter settings.

5.1 Performance and Convergence

Fig. 4 shows the ratio of dropped tasks and the average delay of the proposed DRL-based algorithm and the benchmark methods across episodes (based on the parameters in Table 1). As shown in Fig. 4, the proposed algorithm converges after around 350 episodes, and it achieves a ratio of dropped tasks of 0.02 and an average delay of 0.52 second. This converged performance significantly outperforms those of the benchmark methods. As shown in the figure, the proposed method reduces the ratio of dropped tasks and the average delay by 86.4% – 95.4% and 18.0% – 30.1%, respectively, when compared with the benchmark methods.

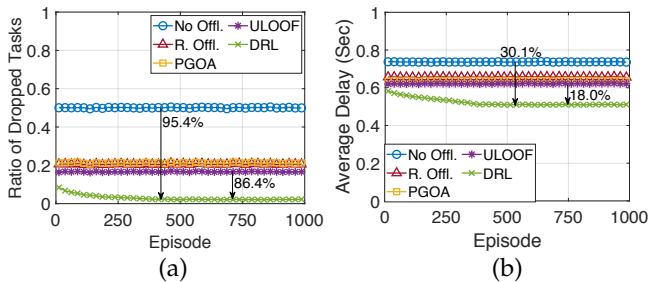


Fig. 4: Performance evaluation across episodes: (a) ratio of dropped tasks; (b) average delay.

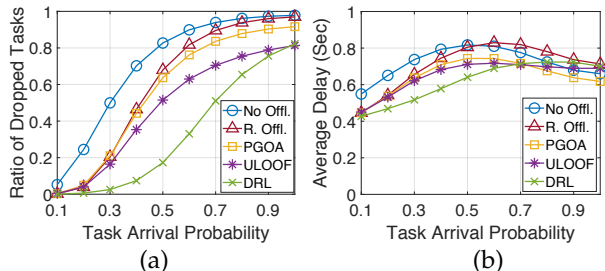


Fig. 5: Performance evaluation under different task arrival probabilities: (a) ratio of dropped tasks; (b) average delay.

5.2 Task Arrival Probability

A larger task arrival probability implies a higher load of the system. As shown in Fig. 5 (a), as the task arrival probability increases, the proposed DRL-based algorithm can always maintain a lower ratio of dropped tasks when compared with the benchmark methods. Specifically, when the task arrival probability is small (i.e., 0.1), most of the methods can achieve a ratio of dropped tasks of around zero. As the task arrival probability increases from 0.1 to 0.5, the ratio of dropped tasks of the proposed algorithm remains less than 0.2, while those of the benchmark methods increase to more than 0.5. In addition, comparing with the benchmark methods, the proposed DRL-based algorithm reduces the ratio of dropped tasks especially when the task arrival probability is moderate (i.e., 0.3 – 0.8), where the reduction of the ratio of dropped tasks is at least 13.3%.

In Fig. 5 (b), as the task arrival probability increases from 0.1 to 0.4, the average delay of our proposed DRL-based algorithm increases by 26.1%, while those of the benchmark methods increase by at least 34.5%. This implies that as the load of the system increases, the average delay of the proposed algorithm increases less dramatically than those of the benchmark methods. As the task arrival probability increases to around 0.6, the average delay of some of the methods decrease, because an increasing number of tasks are dropped and hence are not accounted in the average delay. For the same reason, when the load of the system is high, the proposed algorithm may have a larger average delay than the other methods, as it has less tasks dropped.

5.3 Task Deadline

A smaller deadline implies that the tasks are more delay-sensitive. In Fig. 6 (a), the proposed algorithm always achieves a lower ratio of dropped tasks than the benchmark methods, especially when the deadline is small. When the task deadline is 0.6 second, the proposed algorithm reduces the ratio of dropped tasks by 65.8%–79.3% when compared

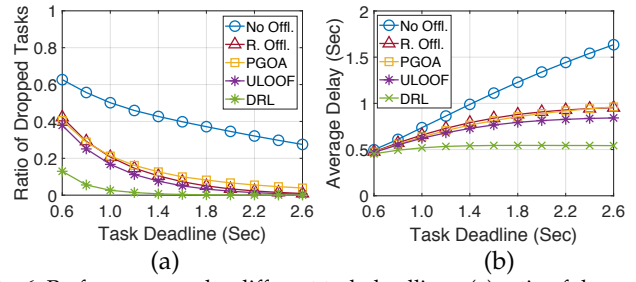


Fig. 6: Performance under different task deadlines: (a) ratio of dropped tasks; (b) average delay.

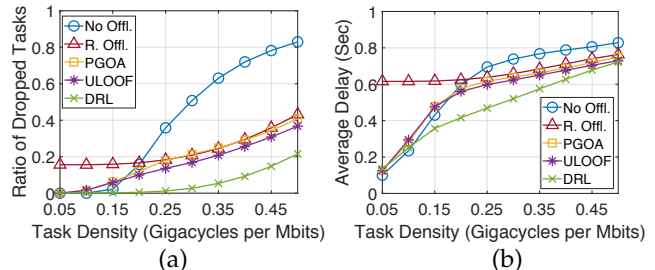


Fig. 7: Performance under different task densities: (a) ratio of dropped tasks; (b) average delay.

with the benchmark methods. As the deadline increases, the ratio of dropped tasks of each method decreases. With the proposed algorithm, the ratio of dropped tasks is less than 0.01 when the deadline is larger than 1.4 seconds. In comparison, the same performance is achieved by ULOOF when the deadline is larger than 2.4 seconds.

In Fig. 6 (b), as the task deadline increases, the average delay of each method increases and gradually converges. This is because when the deadline is larger, the tasks requiring longer processing (and transmission) time can be processed and are accounted in the average delay. When the deadline is large enough, no task is dropped, so further increasing the deadline makes no difference. As shown in Fig. 6 (b), the average delay of the proposed algorithm converges (i.e., achieves a marginal increase of less than 0.05) after the deadline increases to 1.4 seconds, and the converged average delay is around 0.54 second. In comparison, the converged average delay of ULOOF is around 0.84 second, which is 55.6% larger than that of the proposed algorithm, and those of the other methods are larger than 0.96 second. This implies that when the task deadline is large enough, although each method can have a ratio of dropped tasks of around zero, the proposed algorithm outperforms the other methods in terms of reducing the average delay.

5.4 Task Density

A larger task density implies that the computational requirement of each task is larger. As a result, in Fig. 7, as the task density increases, the ratio of dropped tasks and the average delay of each method increase. On the other hand, when the density is small (e.g., smaller than 0.15 Gigacycles per Mbits), the transmission delay dominates the processing delay, so no offloading achieves a lower ratio of dropped tasks and a lower average delay than random offloading. When the density is large (e.g., larger than 0.3 Gigacycles per Mbits), the processing delay dominates the transmission delay, so random offloading achieves a better performance than no offloading.

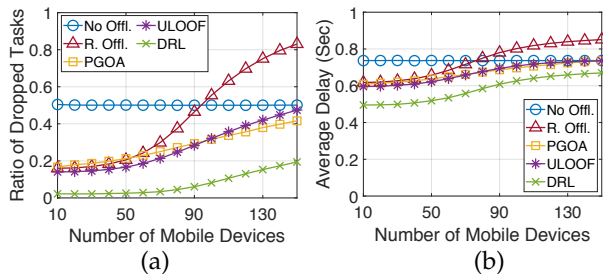


Fig. 8: Performance under different number of mobile devices: (a) ratio of dropped tasks; (b) average delay.

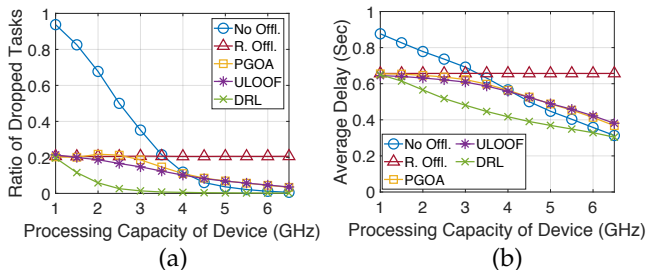


Fig. 9: Performance under different processing capacities of each mobile device: (a) ratio of dropped tasks; (b) average delay.

In Fig. 7, as the task density increases from 0.05 to 0.25 Gigacycles per Mbits, the ratio of dropped tasks and the average delay of the proposed algorithm increase less dramatically than those of the benchmark methods. When the density is 0.25 Gigacycles per Mbits, the proposed algorithm maintains a ratio of dropped tasks of around 0.01 and an average delay of 0.47 second. As the task density further increases to 0.5 Gigacycles per Mbits, although each method achieves a similar average delay, the proposed algorithm can reduce the ratio of dropped tasks by 41.4% – 74.1% when compared with the benchmark methods, because of its proper exploitation of the processing capacities in both the mobile devices and the edge nodes.

5.5 Number of Mobile Devices

A larger number of mobile devices implies potentially a higher load at the edge nodes and hence a worse performance for random offloading. In Fig. 8 (a), the proposed algorithm achieves a lower ratio of dropped tasks than the other methods, especially when the number of mobile devices is large. This is because the proposed algorithm can effectively address the unknown load dynamics at the edge nodes. When the number of mobile devices increases to 80, the proposed algorithm maintains a ratio of dropped tasks of less than 0.05. When it increases to 150, the proposed algorithm achieves a ratio of dropped tasks of 53.4%–76.6% less than the benchmark methods.

In Fig. 8 (b), as the number of mobile devices increases, the average delay of each method (except no offloading) increases due to the potentially increasing load at the edge nodes. Since the proposed algorithm can effectively deal with the unknown edge load dynamics, when the number of mobile devices increases to 150, it achieves an average delay of 9.0% lower than those of PGOA and ULOOF.

5.6 Processing Capacity of Each Mobile Device

As the processing capacity increases, the delay of the tasks processed locally decreases. In Fig. 9 (a), as the processing

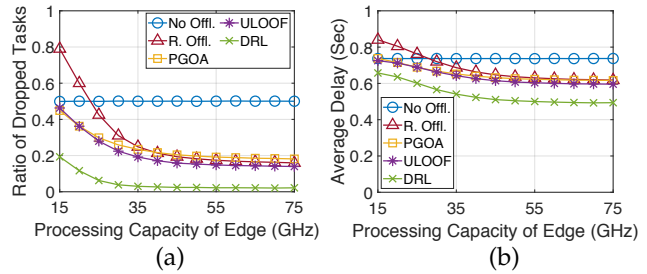


Fig. 10: Performance under different processing capacities of each edge node: (a) ratio of dropped tasks; (b) average delay.

capacity in each device increases, the ratio of dropped tasks of the proposed algorithm decreases more dramatically than those of the benchmark methods. When the processing capacity increases to 3.5 GHz, the ratio of dropped tasks of the proposed algorithm reduces to 0.007, which is 93.9% – 96.5% lower than those of the benchmark methods.

In Fig. 9 (b), as the processing capacity of each device increases, the average delay of the proposed algorithm decreases more dramatically than those of PGOA and ULOOF. When the processing capacity increases to 3.5 GHz, the average delay of the proposed algorithm is 31.4% and 29.4% lower than those of PGOA and ULOOF, respectively. On the other hand, when the processing capacity of each device is large enough, processing a task locally achieves a strictly lower delay than offloading the task to an edge node due to the transmission time required for offloading, and hence no offloading is optimal. Consequently, as the processing capacity increases, the average delay of the proposed algorithm approaches the average delay of no offloading.

5.7 Processing Capacity of Each Edge Node

With a larger processing capacity of each edge node, the average delay of the tasks offloaded is smaller. In Fig. 10, as the processing capacity increases, the ratio of dropped tasks and the average delay of each method (except no offloading) decrease, because an increasing number of tasks are being offloaded. In addition, those values gradually converge. This is because when the processing capacity of each edge node is large enough, further increasing the processing capacity does not reduce the delay of those tasks offloaded due to the limited transmission capacity.

As shown in Fig. 10, the proposed algorithm can reduce the ratio of dropped tasks and the average delay when compared with the benchmark methods. The reduction of the ratio of dropped tasks is especially significant when the processing capacity of each edge node is small. When the processing capacity is 15 GHz, the proposed algorithm reduces the ratio of dropped tasks by at least 57.0% and reduces the average delay by at least 9.4% when compared with the benchmark methods. On the other hand, the converged ratio of dropped tasks and the average delay of the proposed algorithm is at least 84.3% and 17.2% less than those of the benchmark methods, respectively. This is because the proposed algorithm can efficiently exploit the processing capacities in the mobile devices and edge nodes as well as the limited transmission capacity for offloading.

In conclusion, comparing with the benchmark methods, our proposed algorithm achieves a lower ratio of dropped

tasks and a lower average delay under the different parameter settings. The reduction of the ratio of dropped tasks is especially significant when the tasks are delay-sensitive or the load levels at the edge nodes are high (i.e., the task density is large, the number of mobile devices is large, or the processing capacity of each edge node is small).

6 CONCLUSION

In this work, we studied the computational task offloading problem with non-divisible and delay-sensitive tasks in the MEC system and designed a distributed offloading algorithm that enables mobile devices to make their offloading decisions in a decentralized manner. The proposed algorithm can address the unknown load level dynamics at the edge nodes, and it can handle the time-varying system environments (e.g., the arrival of new tasks, the computational requirement of each task). Simulation results showed that when compared with several benchmark methods, our proposed algorithm can reduce the ratio of dropped tasks and average delay. The benefit is especially significant when the tasks are delay-sensitive or the load levels at the edge nodes are high. For future work, it is interesting to enable mobile devices to learn their optimal offloading policies cooperatively through taking advantage of the trained neural networks of other mobile devices. Through the cooperative learning, the training process of the DRL-based algorithm may be accelerated, and the performance may be improved.

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. Surveys & Tuts.*, vol. 19, no. 4, pp. 2322–2358, Aug. 2017.
- [2] F. Bonomi, R. Milito, J. Zhu, and S. Addepalli, "Fog computing and its role in the Internet of things," in *Proc. MCC Workshop on Mobile Cloud Computing (MCC)*, Helsinki, Finland, Aug. 2012.
- [3] P. Porambage, J. Okwuibe, M. Liyanage, M. Ylianttila, and T. Taleb, "Survey on multi-access edge computing for Internet of things realization," *IEEE Commun. Surveys & Tuts.*, vol. 20, no. 4, pp. 2961–2991, Jun. 2018.
- [4] C. Wang, C. Liang, F. R. Yu, Q. Chen, and L. Tang, "Computation offloading and resource allocation in wireless cellular networks with mobile edge computing," *IEEE Trans. Wireless Commun.*, vol. 16, no. 8, pp. 4924–4938, May 2017.
- [5] S. Bi and Y. J. Zhang, "Computation rate maximization for wireless powered mobile-edge computing with binary computation offloading," *IEEE Trans. Wirel. Comm.*, vol. 17, no. 6, pp. 4177–4190, Apr. 2018.
- [6] N. Eshraghi and B. Liang, "Joint offloading decision and resource allocation with uncertain task computing requirement," in *Proc. IEEE INFOCOM*, Paris, France, Apr. 2019.
- [7] X. Lyu, H. Tian, W. Ni, Y. Zhang, P. Zhang, and R. P. Liu, "Energy-efficient admission of delay-sensitive tasks for mobile edge computing," *IEEE Trans. Commun.*, vol. 66, no. 6, pp. 2603–2616, Jun. 2018.
- [8] M. Chen and Y. Hao, "Task offloading for mobile edge computing in software defined ultra-dense network," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 3, pp. 587–597, Mar. 2018.
- [9] K. Poularakis, J. Llorca, A. M. Tulino, I. Taylor, and L. Tassiulas, "Joint service placement and request routing in multi-cell mobile edge computing networks," in *Proc. IEEE INFOCOM*, Paris, France, Apr. 2019.
- [10] X. Lyu, W. Ni, H. Tian, R. P. Liu, X. Wang, G. B. Giannakis, and A. Paulraj, "Distributed online optimization of fog computing for selfish devices with out-of-date information," *IEEE Trans. Wirel. Comm.*, vol. 17, no. 11, pp. 7704–7717, Sep. 2018.
- [11] L. Li, T. Q. Quek, J. Ren, H. H. Yang, Z. Chen, and Y. Zhang, "An incentive-aware job offloading control framework for multi-access edge computing," *IEEE Trans. Mobile Comput.*, 2019 (Early Access).
- [12] H. Shah-Mansouri and V. W.S. Wong, "Hierarchical fog-cloud computing for IoT systems: A computation offloading game," *IEEE Internet of Things J.*, vol. 5, no. 4, pp. 3246–3257, Aug. 2018.
- [13] S. Jošilo and G. Dán, "Wireless and computing resource allocation for selfish computation offloading in edge computing," in *Proc. IEEE INFOCOM*, Paris, France, Apr. 2019.
- [14] L. Yang, H. Zhang, X. Li, H. Ji, and V. 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, Dec. 2018.
- [15] J. L. D. Neto, S.-Y. Yu, D. F. Macedo, M. S. Nogueira, R. Langar, and S. Secci, "ULOOFF: A user level online offloading framework for mobile edge computing," *IEEE Trans. Mobile Comput.*, vol. 17, no. 11, pp. 2660–2674, Nov. 2018.
- [16] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, and G. Ostrovski, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, Feb. 2015.
- [17] L. Huang, S. Bi, and Y. J. Zhang, "Deep reinforcement learning for online computation offloading in wireless powered mobile-edge computing networks," *IEEE Trans. Mobile Comput.*, 2019 (Early Access).
- [18] J. Luo, F. R. Yu, Q. Chen, and L. Tang, "Adaptive video streaming with edge caching and video transcoding over software-defined mobile networks: A deep reinforcement learning approach," *IEEE Trans. Wirel. Comm.*, 2019 (Early Access).
- [19] Y. Liu, H. Yu, S. Xie, and Y. Zhang, "Deep reinforcement learning for offloading and resource allocation in vehicle edge computing and networks," *IEEE Trans. Veh. Technol.*, vol. 68, no. 11, pp. 11 158–11 168, Nov. 2019.
- [20] 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. Comm.*, vol. 18, no. 11, pp. 5141–5152, Nov. 2019.
- [21] J. Liu, Y. Mao, J. Zhang, and K. B. Letaief, "Delay-optimal computation task scheduling for mobile-edge computing systems," in *Proc. IEEE Int'l Symposium on Information Theory (ISIT)*, Barcelona, Spain, Jul. 2016.
- [22] T. Ouyang, Z. Zhou, and X. Chen, "Follow me at the edge: Mobility-aware dynamic service placement for mobile edge computing," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 10, pp. 2333–2345, Oct. 2018.
- [23] S. Jošilo and G. Dán, "Selfish decentralized computation offloading for mobile cloud computing in dense wireless networks," *IEEE Trans. Mobile Comput.*, vol. 18, no. 1, pp. 207–220, Jan. 2019.
- [24] 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.
- [25] A. Demers, S. Keshav, and S. Shenker, "Analysis and simulation of a fair queueing algorithm," *ACM SIGCOMM Computer Communication Review*, vol. 19, no. 4, pp. 1–12, Sep. 1989.
- [26] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1997.
- [27] F. A. Gers, J. Schmidhuber, and F. Cummins, "Learning to forget: Continual prediction with LSTM," in *Proc. Int'l Conf. on Artificial Neural Networks (ICANN)*, Edinburgh, UK, Sep. 1999.
- [28] Z. Wang, T. Schaul, M. Hessel, H. van Hasselt, M. Lanctot, and N. de Freitas, "Dueling network architectures for deep reinforcement learning," in *Proc. Int'l Conf. on Machine Learning (ICML)*, New York City, NY, Jun. 2016.
- [29] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing Atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.
- [30] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.
- [31] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [32] H. van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double Q-learning," in *Proc. AAAI Conf. on Artificial Intelligence*, Phoenix, AZ, May 2016.
- [33] Speedtest Intelligence, "Speedtest market reports: Canada average mobile upload speed based on Q2-Q3 2019 data," <https://www.speedtest.net/reports/canada/>, accessed on Mar. 27, 2019.