

Edge-Aided Computing and Transmission Scheduling for LTE-U-Enabled IoT

Hongli He, Hanguan Shan, Aiping Huang, Qiang Ye, Weihua Zhuang

Abstract

To facilitate the deployment of private industrial Internet-of-Things (IoT), applying long-term-evolution (LTE) over unlicensed spectrum (LTE-U) is a promising technology, which can deal with the licensed spectrum scarcity problem and the stringent quality-of-service (QoS) requirement via centralized control. In this paper, we investigate the computing offloading problem for LTE-U-enabled IoT, where computing tasks on an IoT device are either executed locally or offloaded to the edge server on an LTE-U base station. Considering a constrained edge computing cost (e.g., operation power consumption) for offloaded tasks, the task scheduling problem is formulated as a constrained Markov decision process (CMDP) to maximize the long-term average reward, which integrates both task completion profit and task completion delay. In order to address the uncertainty of task arrivals and channel availability, a constrained deep Q-learning-based task scheduling algorithm with provable convergence is proposed, where an adaptive reward function can appropriately bound the average edge computing cost. Extensive simulation results show that the proposed scheme considerably enhances the system performance.

Index Terms

Mobile edge computing, task offloading, Internet-of-Things, LTE over unlicensed, constrained deep reinforcement learning.

I. INTRODUCTION

Recently, there has been a soaring proliferation of smart devices with ubiquitous interconnections for providing Internet-of-Things (IoT) services [2]. By enabling the osmotic convergence of transmission, computing, and storage on IoT devices, a diversity of promising applications have emerged, encompassing smart homing, intelligent connected vehicles, and wearable IoT systems [3]. While many IoT applications can be supported via cellular networks, there are certain entities that prefer their own private networks for reasons such as data security and/or network controllability [4]. Without having licensed radio spectrum resources, these entities generally resort to operate their private networks using unlicensed spectrum. One solution is to

This work was presented in part at IEEE Globecom 2018 [1].

1
2
3 use non-3GPP air interface technologies, such as low-power wide-area (LPWA)-based systems
4 over the unlicensed spectrum [5]. Most existing LPWA-based systems such as long range
5 (LoRa) over the unlicensed spectrum are usually simplified for cost-saving purposes, at the
6 price of service performance [6]. On the other hand, long-term evolution (LTE) over unlicensed
7 spectrum (LTE-U) provides an alternative solution that can share the cellular IoT ecosystem with
8 unlicensed spectrum resources and satisfy the enterprise customers' requirements for a dedicated
9 IoT network. Technically, many of the advanced wireless techniques of LTE-U inherited from
10 cellular systems are proven to offer better performance for IoT connections than traditional IoT
11 technologies over the unlicensed spectrum, and overcome performance limitations of existing
12 LPWA systems [7]. For better promoting the performance of LTE-U network in IoT services, a lot
13 of organizations, especially Multefire, have proposed some dedicated modifications, for example,
14 supporting eMTC (with bandwidth of 1.4 MHz) in unlicensed bands, supporting narrow band
15 (NB)-IoT (with bandwidth of 200kHz) in unlicensed bands, and being compatible with lower
16 spectrum bands (sub-1GHz). Currently, LTE-U-based technology has been deployed in some
17 practical IoT systems such as Shanghai Yangshan Phase VI Port, where the LTE-U network can
18 support the automated guided vehicle controlling, wireless close-circuit Television (CCTV), and
19 remote crane monitoring system in an area with length of 2350 meters and seven berths [8].

20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Meanwhile, the constrained computational capability and power capacity of IoT devices pose a significant challenge for the quality-of-service (QoS) guarantee of the emerging applications, especially for those with intensive computation, e.g., image processing in multimedia IoT, although local computing at the device can significantly reduce the latency for transmission and the privacy leakage [9]. One simple way to address this issue is to establish a sustainable tunnel between the resource-limited devices and the cloud computing system. However, cloud servers are usually located far away from the end-devices, leading to high transmission delay over backhaul link and excessive stress of traffic load on the core networks [10]. To tackle these challenges, a mobile edge computing (MEC) paradigm emerges via shifting some tasks from the devices to vicinal computation capable entities, i.e., the edge of wireless networks. The advent of MEC facilitates the extension of cloud computing services to the edge within the radio access network of the end-devices, and provides a swift computing and mobility-aware support [11].

However, employing MEC for LTE-U-enabled IoT faces significant challenges to satisfy the high reliability and low latency requirements. Specifically, the task forwarding from an IoT device to the edge server can be interrupted by other transmissions from the legacy networks

1
2
3 (e.g. a Wi-Fi network) operating over the unlicensed bands, leading to a performance degradation.
4 Further, the carrier-sensing-based channel access in LTE-U networks causes extra waiting time
5 of detecting channel availability in the reverse transmission from the edge server to the IoT
6 device, while the delay for the reverse transmission over licensed spectrum is not considered
7 in existing studies due to small data size of task execution results. Therefore, dynamics of the
8 unlicensed channel availability lead to long delay and unreliability in data transmission of the
9 MEC system, and make the delay analysis intractable.
10
11
12
13

14 In this paper, we study an edge-aided task scheduling problem in an LTE-U-enabled network,
15 where an LTE-U base station (BS) helps IoT devices to decide whether tasks should be locally
16 computed or be offloaded to the edge server on the LTE-U BS. Specifically, by taking account
17 of the dynamic queueing, transmission, and computing of the associated tasks, we propose a
18 state-wise model for the evolution of the MEC system and the availability of wireless channels.
19 Correspondingly, the profit of task completion and the joint computing and transmission delay
20 are integrated for evaluating the effect of different task scheduling policies, under a constraint of
21 the edge computing cost (e.g., energy consumption for offloaded tasks). The main contributions
22 of this paper are threefolded:
23
24
25
26
27
28
29

- 30 • The applicability of the constrained Markov decision process (CMDP) framework is ex-
31 plored for the task scheduling in an LTE-U-enabled IoT network for attaining a high reward
32 related to the profit and delay of task completion, with a constraint of the average edge
33 computing cost. The random generation, transmission, and computation of the tasks, the
34 status of the edge server, and the dynamics of the unlicensed channel availability are captured
35 in the proposed stochastic optimization-based task scheduling framework;
36
37
38
39
40
- 41 • With no priori knowledge about channel status and task arrivals, we propose a constrained
42 deep Q-learning (CDQL) solution for the task scheduling problem. Motivated by applying
43 the Lagrange duality in solving CMDP problem with available environment information,
44 we propose, for the first time, a subgradient-based reinforcement learning algorithm to
45 solve the cost-constrained task scheduling, which is transferred to a standard reinforcement
46 learning problem without cost constraint when the optimal Lagrange dual multiplier is given.
47 Specifically, the value of the Lagrange dual multiplier is proposed to be updated based on
48 the counted average edge computing cost when the environment information is unavailable,
49 and the convergence of the subgradient-based learning algorithm is proved;
50
51
52
53
54
55
- 56 • For implementing the proposed CDQL algorithm, we modify the traditional deep Q-learning
57
58
59
60

(DQL) architecture including the deep neural network (DNN) structure for estimating the Q-values, the preprocessing of the training samples, and the pretraining with adaptive learning rates to accelerate the learning for Q values.

The remainder of the paper is organized as follows: We discuss the related work in Section II. System model is described in Section III. In Section IV, we formulate the edge-aided task scheduling problem as a CMDP framework and present the CDQL algorithm in Section V. The customized modification for the realization of our proposed algorithm is presented in Section VI and simulation results are given in Section VII, followed by conclusions in Section VIII.

II. RELATED WORK

There have been a collection of research works studying how to guarantee the performance of IoT over the unlicensed spectrum. For achieving extended coverage, low power consumption, and massive connectivity over unlicensed spectrum, NB-IoT-U, a cellular industrial IoT system over unlicensed spectrum, is proposed by expanding NB-IoT technology in sub-1-GHz [7]. In [12], a distributed channel sensing and resource allocation scheme is proposed to achieve a high resource reuse ratio for massive IoT connections via appropriately mitigating intra/inter-network interference over the unlicensed spectrum. In order to extend the lifetime of IoT devices, energy harvesting is enabled and analyzed for an IoT network over unlicensed spectrum in [13]. In [14], a cellular-user-aided relay scheme is proposed to bridge the communication between an IoT device and its destined BS via machine-to-machine communication over the unlicensed spectrum, which reduces the energy consumption and supports more connected devices. These works have addressed many communication-related issues for the LTE-U-enabled IoT, which builds a robust radio access network for supporting services in terms of coverage, connectivity, interference coordination, and energy consumption. However, accommodating computing-intensive services in LTE-U-based networks still faces significant challenges due to the unreliable communication links over the unlicensed spectrum and the complex orchestration for both transmission and computing resources [15].

While the MEC in LTE-U-enabled IoT is under investigated, applying MEC for licensed IoT has been exploited extensively. The existing researches mainly use two kinds of methods to design the task scheduling strategy for edge computing. The first one is Lyapunov-based method. In [16], by integrating the queue length of edge servers, Chen *et al.* propose a dynamic computing offloading scheme to achieve the tradeoff between offloading cost and performance based on Lya-

punov optimization. In [17], considering the uncertainty and dynamics of wireless channel state and task arrival process, and the large scale of solution space, an energy-efficient task offloading scheme is proposed by leveraging the Lyapunov drift. To relieve the stringent requirement of tasks' feedback, a scheduling algorithm tolerant to out-of-date network knowledge is proposed in [18] and can achieve asymptotical optimality with only partial information via applying both the perturbed Lyapunov function and the knapsack problem formulation. The queueing-theory-based and Lyapunov-drift-based model in these papers can maintain the theoretical stability of the MEC system when the task generation process is dynamic. However, such models cannot integrate some non-queue-based information, e.g., the channel access state of the LTE-U network when employing listen-before-talk (LBT)-based channel access schemes, and neglect the exploration for the network statistics in a large time scale (e.g., task generation). Therefore, such algorithms usually sacrifice their performance for the system robustness. The second one is the reinforcement learning method. In [19], a reinforcement-learning-based scheduling scheme is proposed for the virtual machine assignment and task offloading in a space-air-ground integrated network. In [20], considering the stochastic vehicle traffic, dynamic computation requests, and time-varying communication conditions, a semi-Markov process is formulated and a deep reinforcement learning method is proposed to obtain the optimal policies of computation offloading and resource allocation. In [21], a multi-user multi-edge-node computation offloading problem is formulated as an unknown payoff game and solved by the distributed reinforcement learning method. In these works, the reinforcement learning methods are mostly based on the Q-value framework, and cannot be directly applied in the scenario with an edge computing cost constraint. In the following, we present a comprehensive CDQL-based task scheduling framework for an LTE-U-enabled IoT with MEC to facilitate the state-wise task offloading decision while guaranteeing the edge computing cost constrained within the given cost budget.

III. SYSTEM MODEL

A. Network Scenario

We consider an MEC-aided LTE-U-enabled IoT as shown in Fig. 1, consisting of one LTE-U BS, and K LTE-U-enabled IoT devices. Let the indexes of the K devices be denoted by the elements in $\mathcal{K} = \{1, 2, \dots, K\}$ and the index of the BS be 0. Transmission and computing functions are both integrated in the LTE-U BS operating over unlicensed spectrum. A channel with bandwidth B is assigned for the MEC service supported by an edge server on the LTE-U

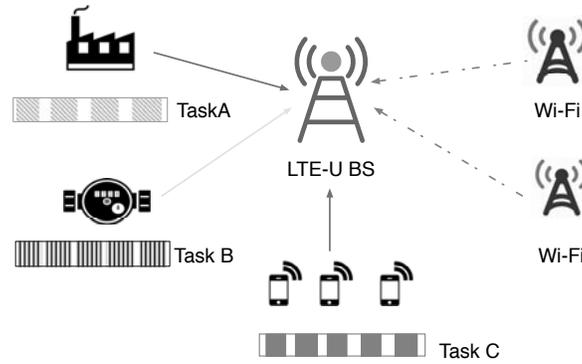


Fig. 1: The deployment scenario.

BS. The system evolves over time in a slotted structure indexed by t (> 0) and the length of each slot is τ . Because the IoT devices are normally deployed for some dedicated services [22], we consider that each device has only one application and needs to keep performing one kind of corresponding task [23]. Due to the stringent computing capability of the devices and the possible low latency requirements, some tasks should be offloaded to the edge server on the LTE-U BS which is empowered with much richer computing resources.

The structure of tasks generated from the same device has a consistent format that can be processed by corresponding devices. The task structure at device k is denoted by $M_k \triangleq (L_k, X_k, f_k)$, where L_k denotes the size of the task input (in bit), X_k is the required computation intensity (in CPU cycles per bit), and f_k is the profit for task completion [24]. These statistics can be obtained using the method proposed in [25]. A linear dependency is considered for two consecutive tasks of the same device, i.e., a task of device k ($\in \mathcal{K}$) can only be performed after the completion of its previous task. In a practical scenario, an urgent task with high priority can preempt the unserved low-priority tasks, which changes the original serving order. However, such special cases do not affect the system modeling under consideration in this work because the computing of the reordered tasks leads to the same reward (as specified in Section IV) when the scheduling policy is fixed. By setting the duration of each time slot as a small value, it is assumed that at most one task is generated at a device within each time slot. The task arrival probability at device k in slot t is denoted by $p_k(t)$. The computing capacity (in the unit of cycles per second) of device k is represented by C_k ($k \in \mathcal{K}$) while the computing capacity of the BS is C_0 ($C_0 \gg C_k, \forall k \in \mathcal{K}$).

Since an LTE-U network shares the unlicensed spectrum with other legacy networks, e.g.,

1
2
3 Wi-Fi, for fair coexisting, an LBT mechanism is applied for the LTE-U network [26]. By taking
4 into account both flexibility and simplicity of the channel access, a category 4 scheme (LBT
5 with random back-off and a contention window of variable size) defined by 3GPP is adopted
6 in the LTE-U BS for the unlicensed channel access [27]. The listening (i.e., channel sensing)
7 of the LTE-U network is executed by the BS and the channel-related state is represented by
8 tuple $s_C(t) = (a(t), b(t), m(t))$, where $a(t)$ is the phase indicator, $b(t)$ is the phase counter, and
9 $m(t)$ is the backoff stage. When $a(t) = 0$ (or 1), the BS is under the sensing (or transmission)
10 phase. Meanwhile, $b(t)$ is a backoff counter with different indications under different phases.
11 When the BS intends to access the channel, the backoff stage $m(t)$ is set as 0, and the backoff
12 counter $b(t)$ is randomly initialized between 0 and $W_0 - 1$. Then the BS senses the channel at the
13 beginning of each slot. The backoff counter $b(t)$ is decreased by 1 at each time slot under idle
14 channel, while it keeps frozen if the channel is busy. If there is no transmission collision with
15 other networks over the unlicensed channel after the backoff counter $b(t)$ is reduced to 0, the
16 BS begins to access the unlicensed channel for scheduling the devices to do their task-related
17 transmission. Under the collision case, the backoff stage $m(t)$ is increased by 1, and the backoff
18 window is doubled, i.e., $W_{m(t)} = 2W_{m(t)-1}$. The maximum backoff stage is M . Then the backoff
19 counter $b(t)$ is randomly reset between 0 and $W_{m(t)} - 1$ for the subsequent channel sensing.
20
21
22
23
24
25
26
27
28
29
30
31

32 After successfully receiving the uplink grant signal, the associated IoT devices immediately
33 upload their task data and equally share the available frequency bands. The uplink transmission
34 rate of device k at slot t is given by
35
36

$$37 \quad r_k(t) = \frac{B}{J(t)} \log \left(1 + \frac{P_k |h_k(t)|^2}{\sigma^2} \right) \quad (1)$$

38 where $J(t)$ is number of devices with transmission demand in slot t , P_k is the fixed transmission
39 power of device k , $h_k(t)$ is the channel gain from device k to the BS, including the path loss
40 and Rayleigh fading, and σ^2 is the received noise power. At the end of the transmission phase,
41 $a(t)$ and $b(t)$ are reset to 0 and W_0 , respectively, for the next round of channel sensing. The
42 Wi-Fi system works as one external environmental element. Specifically, as the number of Wi-Fi
43 users increases, the average channel quality of the LTE-U decreases, which leads to a longer
44 transmission delay. Hence, the tasks are more likely to be processed locally. That is, the activity
45 of Wi-Fi devices affects the evolution of channel availability, hence the transmission queue size,
46 and finally the offloading decision.
47
48
49
50
51
52
53
54
55
56
57
58
59
60

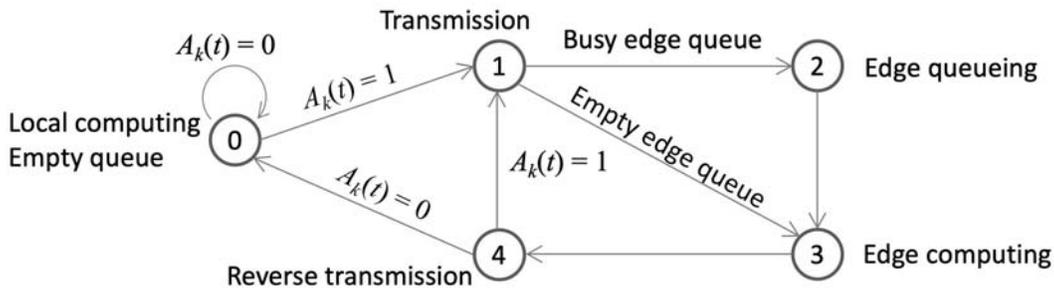


Fig. 2: The task processing flow.

B. Offloading Policy and Device Status

When a task becomes the head of the task queue in device k at the beginning of time slot t , it is scheduled by the BS for locally processing or offloading. The locally processing at device k is denoted by action $A_k(t) = 0$ while offloading decision is denoted by $A_k(t) = 1$. The decision epochs for task processing are discussed in Section III-C. The device status is highly dependent on its task that is being processed, referred to as in-service task. The status of all devices is captured by vector $\mathbf{z}(t) = (z_1(t), z_2(t), \dots, z_k(t), \dots, z_K(t))$, $z_k(t) \in \{0, 1, 2, 3, 4\}$, $k \in \mathcal{K}$. Specifically, for in-service task of device k in slot t , $z_k(t) = 0$ means that the task is being locally processed (including the case that the task queue is empty); $z_k(t) = 1$ means that the device is uploading the input data of the task to the BS; $z_k(t) = 2$ denotes that the in-service task of device k is queueing at the edge server; $z_k(t) = 3$ denotes that the task is being processed at the edge server; $z_k(t) = 4$ indicates that the task is accomplished by the edge server and is waiting for returning task execution result from the edge server. As shown in Fig. 2, the evolution for the device status depends on the availability of the unlicensed channel and computing resources, as discussed in the following subsections.

C. Task Backlog

Each device is equipped with a task buffer of capacity \bar{B}_k , containing all its task backlog. The number of unaccomplished tasks in the buffer of device k at the beginning of time slot t is denoted by $B_k(t)$. The generation of new tasks and the processing of in-service task jointly decide backlog update, expressed as

$$B_k(t+1) = \min \{B'_k(t) + Y_k(t), \bar{B}_k\} \quad (2)$$

where $Y_k(t)$ is a Bernoulli-distributed random variable with parameter $p_k(t)$, indicating whether

in slot t a new task is generated at device k , and $B'_k(t)$ denotes the updated backlog size without new task generation, derived by

$$B'_k(t) = \begin{cases} \max\left\{B_k(t) - \frac{C_k\tau}{L_k X_k}, \lfloor B_k(t) \rfloor^-\right\}, & \text{if } z_k(t) = 0 \\ B_k(t) - 1, & \text{if } z_k(t) = 4 \text{ and } a(t) = 1 \\ B_k(t), & \text{otherwise} \end{cases} \quad (3)$$

where $\lfloor \cdot \rfloor^-$ is the revised floor function mapping to the greatest integer less than itself. The first subequation in (3) represents the case that a task is being locally computed at device k , and its backlog size is reduced by the normalized fraction. The second subequation indicates that an in-service task of device k has been completed at the edge server, and the backlog is updated when the task execution result is returned from the edge server once the channel is available.

The decision epochs for device k can be classified into two categories. One is referred to as time slots with a local task completion or a return of the task output from the edge, while the local task backlog is not empty, which is expressed by $\mathcal{H}_k = \{t | t \in \mathcal{T}_k^0, B_k(t+1) > 0\}$, $\mathcal{T}_k^0 = \{t | B'_k(t) \in \mathbf{Z}, B'_k(t) < B_k(t)\}$; The other is the case that a new task is generated during slot t at device k when its task backlog is empty at the slot beginning, which is denoted by $\mathcal{E}_k = \{t | B_k(t) = 0, Y_k(t) > 0\}$. Therefore, the overall decision epochs for device k are denoted as $\mathcal{E}_k = \mathcal{H}_k \cup \mathcal{E}_k$.

D. Local Transmission Queue of Device

After a task of device k is scheduled to be offloaded, device k pushes the task input data into its transmission queue for uploading to the edge server. We normalize the transmission queue length of device k at the beginning of slot t with regard to L_k , which is denoted by $F_k(t)$. Because the tasks are scheduled sequentially, there is at most one task in the transmission queue, i.e., $F_k(t) \in [0, 1]$. The evolution of the transmission queue length is given by

$$F_k(t+1) = \begin{cases} \max\left\{F_k(t) - \frac{r_k(t)\tau}{L_k}, 0\right\}, & \text{if } z_k(t) = 1 \text{ and } a(t) = 1 \\ F_k(t) + 1, & \text{if } A_k(t) = 1 \\ F_k(t), & \text{otherwise.} \end{cases} \quad (4)$$

The first subequation in (4) means that the transmission queue length of device k is decreased due to the uplink transmission of its task input data when the LTE-U BS is in the transmission

1
2
3 phase. The second subequation indicates that the transmission queue length is added by 1 if
4 a new task is to be offloaded to the edge server ($A_k(t) = 1$). The epochs of transmission
5 completion of device k are denoted as $\mathcal{T}_k^1 = \{t | F_k(t+1) = 0, F_k(t) > 0\}$ and we have the
6 overall transmission completion epochs for any device as $\mathcal{T}^1 = \bigcup_{k \in \mathcal{K}} \mathcal{T}_k^1$.
7
8
9

10 *E. Edge Computing Queue at Edge Server*

11 After the uplink transmission for the task input of device k is finished, the task is pushed
12 into the edge computing queue. A first-come-first-serve edge computing queue is considered. To
13 guarantee that the tasks are served in an appropriate sequence, an edge computing priority $\mathcal{W}_k(t)$
14 ($k \in \mathcal{K}$) is used for the in-service task from device k at time t . Once any new task finishes its
15 uploading for its input data and enters the edge computing queue, the priority of all tasks at
16 the edge (including the new task) are increased by one, guaranteeing that the earlier tasks are
17 allocated with higher priority. The update of the edge computing priority is denoted by
18
19
20
21
22
23

$$24 \quad \bar{\mathcal{W}}_k(t) = \begin{cases} \mathcal{W}_k(t) + 1, & \text{if } t \in \mathcal{T}^1 \text{ and } k \in \mathcal{C}(t) \\ \mathcal{W}_k(t), & \text{otherwise} \end{cases} \quad (5)$$

25 where $\mathcal{C}(t)$ is the set containing devices with tasks in the edge computing queue, i.e., $\mathcal{C}(t) =$
26 $\{k | \mathcal{W}_k(t) > 0 \text{ or } t \in \mathcal{T}_k^1\}$. Moreover, for integrating the computing process of tasks at the edge,
27 the decimal part of $\mathcal{W}_k(t)$ is interpreted as the unfinished task fraction of device k 's in-service
28 task at the edge server, which is normalized by $L_k X_k$. The computing-related update for the
29 priority is described as
30
31
32
33
34
35
36
37

$$38 \quad \hat{\mathcal{W}}_k(t) = \begin{cases} \max \left\{ \bar{\mathcal{W}}_k(t) - \frac{C_0 \tau}{L_k X_k}, [\bar{\mathcal{W}}_k(t)]^- \right\}, & \text{if } z_k(t) = 3 \\ \bar{\mathcal{W}}_k(t), & \text{otherwise.} \end{cases} \quad (6)$$

39 In (6), the first subequation means that, if a task is being processed at the edge server ($z_k(t) = 3$),
40 the decimal part of its priority value is reduced by the normalized processed fraction. Once a
41 task releases the edge computing resources due to the task accomplishment, its priority value is
42 reset to zero, which is denoted by
43
44
45
46
47
48
49

$$50 \quad \mathcal{W}_k(t+1) = \begin{cases} 0, & \text{if } \hat{\mathcal{W}}_k(t) \in \mathbf{Z} \text{ and } z_k(t) = 3 \\ \hat{\mathcal{W}}_k(t), & \text{otherwise.} \end{cases} \quad (7)$$

51 The task completion epochs of edge computing for device k are denoted by set $\mathcal{T}_k^2 = \{t | \hat{\mathcal{W}}_k(t) \in$
52 $\mathbf{Z}, z_k(t) = 3\}$. Since the data size of task execution result is normally much smaller than the
53
54
55
56
57
58
59
60

input data size, the transmission of the output data for all tasks is assumed to take only one time slot [28].

F. Delay Counter

Suppose that task i of device k is generated at T_i^A , scheduled at T_i^S , and completed at T_i^E . The delay of this task is calculated as $d_i = T_i^E - T_i^A$. For recording the delay of each task, a collection of counters should be maintained in each device, which significantly increases the modeling complexity. Fortunately, we can exploit the correlation between two successive tasks in the same device to represent the delay in a more effective way. Suppose that we have only one delay counter D_k representing the number of slots between the generation epoch of the in-service task of device k and the current slot. For task $i + 1$ generated at T_{i+1}^A , scheduled at T_{i+1}^S , and completed at T_{i+1}^E , its scheduling delay of task $i + 1$ is given by

$$d_{i+1} = T_{i+1}^E - T_{i+1}^A = T_{i+1}^E - T_{i+1}^S + T_{i+1}^S - T_{i+1}^A. \quad (8)$$

If task i has been completed when task $i + 1$ is generated, the latter enters an empty task backlog and the delay counter of device k is directly used to record the delay of task $i + 1$. If the task buffer is not empty when task $i + 1$ is generated, the delay counter is being used by the former task i . The history scheduling epochs and ending epochs before task $i + 1$ are respectively denoted by $\mathbf{T}_i^S = (T_{i'}^S)_{i' \leq i}$ and $\mathbf{T}_i^E = (T_{i'}^E)_{i' \leq i}$. Actually, considering that the ending epoch of task i is actually the scheduling epoch of task $i + 1$, i.e., $T_{i+1}^S = T_i^E$, we have

$$d_{i+1} = T_{i+1}^E - T_{i+1}^S + T_i^E - (T_i^A + \Delta_i) = T_{i+1}^E - T_{i+1}^S + d_i - \Delta_i \quad (9)$$

where $\Delta_i = T_{i+1}^A - T_i^A$ is the interval between the generations of tasks i and $i + 1$. Considering that the task arrivals of each device in a period of interest can be modeled as a stationary stochastic process, the expectation of Δ_i is denoted by $\bar{\Delta}$. We calculate the expected delay of task $i + 1$ under the condition of the real measured scheduling and ending epochs history as

$$\mathbb{E}(d_{i+1} | \mathbf{T}_{i+1}^S, \mathbf{T}_{i+1}^E) = \underbrace{T_{i+1}^E - T_{i+1}^S}_{\text{after scheduling}} + \underbrace{\mathbb{E}(d_i | \mathbf{T}_i^S, \mathbf{T}_i^E) - \bar{\Delta}}_{\text{before scheduling}}. \quad (10)$$

From (10), the expected delay of a task is divided by the time before scheduling and that after scheduling. The expected waiting time before scheduling is calculated based on the expectation of previous task's delay ($\mathbb{E}(d_i | \mathbf{T}_i^S, \mathbf{T}_i^E)$) and the average interval of two successive task arrivals ($\bar{\Delta}$). Therefore, we only set one delay counter for each IoT device for counting the expected

1
2
3 delay of the in-service task conditioned on the epochs of task scheduling and completion. Upon
4 completing the current task and scheduling the next task, the delay counter is reset to the expected
5 waiting time ($\mathbb{E}(d_i | \mathbf{T}_i^S, \mathbf{T}_i^E) - \bar{\Delta}$) before scheduling for the next task. Thus, we obtain the average
6 delay of all tasks, sequentially, by only using the delay counter information of the head-of-line
7 task. The evolution of the delay counter is given by
8
9

$$10 \quad D_k(t+1) = \begin{cases} D_k(t) - \bar{\Delta}, & \text{if } t \in \mathcal{H}_k \\ 0, & \text{if } t \in \mathcal{E}_k \\ D_k(t) + 1, & \text{if } B_k(t) > 0 \text{ and } t \notin \mathcal{E}_k \\ D_k(t), & \text{otherwise.} \end{cases} \quad (11)$$

11
12
13
14
15
16
17
18
19
20 In (11) the first subequation indicates that a new task of device k is scheduled at the end of
21 slot t and its delay counter is initialized as $D_k(t) - \bar{\Delta}$; the second subequation indicates that
22 when the task buffer is empty after a task completion at the end of slot t , the delay counter
23 is reset to 0; the third subequation indicates that, when the task is being scheduled, the delay
24 is increased by one after each time slot. It is noted that the delay counter can be negative due
25 to the estimation error between the true sampling interval of the two successive tasks and its
26 expectation. However, this potential negative value can be balanced in the long-term evolution.
27 Different from directly estimating the delay information based on LTE-U/Wi-Fi coexistence
28 analysis such as in [29, 30] which needs the specific information of Wi-Fi network, setting the
29 delay counter at each device can provide decision maker with extra real-time information and
30 make the task scheduling strategy more intelligently without the coordination between LTE-U
31 and W-Fi.
32
33
34
35
36
37
38
39
40

41 IV. PROBLEM FORMULATION

42
43 To capture the dynamic evolution of task-related states, we leverage a Markov decision
44 process (MDP) framework to model the complex interactions among the processes of trans-
45 mission, computing, and queueing in the MEC system. Basically, an MDP is defined by a
46 tuple of state space \mathcal{S} , decision space \mathcal{A} , state transition probability function $P := \mathcal{S} \times \mathcal{A} \times$
47 $\mathcal{S} \rightarrow \mathbb{R}$, and reward function $R := \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$. The system state is defined as $\mathbf{s}(t) =$
48 $(\mathbf{z}(t), \mathbf{B}(t), \mathbf{F}(t), \mathbf{W}(t), \mathbf{D}(t), a(t), b(t), m(t))$, where $\mathbf{B}(t)$, $\mathbf{F}(t)$, $\mathbf{W}(t)$, and $\mathbf{D}(t)$ are the vec-
49 torized values of the task backlog sizes, transmission queue sizes, edge computing priorities,
50 and delay counters, respectively, of all devices at time t , and their values can be obtained from
51
52
53
54
55
56
57
58
59
60

the derivations in Section III. In each generalized decision epoch, denoted by $t \in \mathcal{E} = \bigcup_k \mathcal{E}_k$, an action $\mathbf{A}(t) = [A_k(t)]_{k \in \mathcal{K}}$ is taken. The available action space is dependent on the current state and is expressed as

$$\mathcal{A}(t) = \prod_{\otimes} \mathcal{A}_k(t) \quad (12)$$

where \prod_{\otimes} is the cumulative Cartesian product operation, and $\mathcal{A}_k(t)$ is the sub-space containing the available actions for device k . Specifically, we have

$$\mathcal{A}_k(t) = \begin{cases} \{0, 1\} & \text{if } t \in \mathcal{E}_k \\ \{0\} & \text{otherwise} \end{cases} \quad (13)$$

which indicates that the decision space of device k includes both local computing ($A_k(t) = 0$) and offloading ($A_k(t) = 1$) when the current slot is a decision epoch of the device, otherwise, it is just set as 0 to keep the notation consistent.

The state transitions consist of three steps: The first one is to update $\mathbf{B}(t)$, $\mathbf{F}(t)$, $\mathbf{W}(t)$, and $\mathbf{D}(t)$; The second one is to update the channel status based on the transition probability of different states, which depends on the external environment and the LBT setting of the LTE-U BS [31]; The third one is to update the task status, given by

$$z_k(t+1) = \begin{cases} A_k(t), & \text{if } t \in \mathcal{E}_k \\ 2, & \text{if } z_k(t) = 1 \text{ and } t \in \mathcal{T}_k^1 \text{ and } k \neq \arg \max_{k \in \mathcal{K}} \mathcal{W}_k(t+1) \\ 3, & \text{if } z_k(t) = 1 \text{ and } t \in \mathcal{T}_k^1 \text{ and } k = \arg \max_{k \in \mathcal{K}} \mathcal{W}_k(t+1), \text{ or} \\ & \text{if } z_k(t) = 2 \text{ and } t \in \mathcal{T}^2 \text{ and } k = \arg \max_{k \in \mathcal{K}} \mathcal{W}_k(t+1) \\ 4, & \text{if } z_k(t) = 3 \text{ and } t \in \mathcal{T}_k^2 \\ z_k(t), & \text{otherwise.} \end{cases} \quad (14)$$

In (14), the first case represents that, in a decision epoch, the task status changes into the local computing status if $A_k(t) = 0$ or the transmission status if $A_k(t) = 1$; The second case indicates that the task transits into the edge queueing status after it finishes uploading the task input data and the edge server is occupied; In the third case, if the edge computing queue is empty, the task directly begins to be served by the edge server after its transmission of the task input, or the task status transits from edge queueing status to the edge computing status when other tasks with higher priority have all been completed; The fourth case means that a completed task of

device k begins to wait at the BS for returning its output data back to device k when the channel is accessible.

After each task is completed under state $\mathbf{s}(t)$, a reward $R_k(\mathbf{s}(t))$ is generated for device k , including both the positive profit $(1 - \eta)f_k$ and the negative delay-related impact $-\eta D_k$, given by

$$R_k(\mathbf{s}(t)) = \begin{cases} (1 - \eta)f_k - \eta D_k(t), & \text{if } t \in \mathcal{T}_k^0 \\ 0, & \text{otherwise} \end{cases} \quad (15)$$

where $\eta \in [0, 1]$ is a weight parameter to balance the tradeoff between the profit and delay of the task completion. The reward function can be customized in different forms for the worst-case delay guarantee. For example, in [32], a hard delay deadline is set for offloading or local execution of each task. If a task computing delay exceeds the deadline, the task would be dropped, and a large dropping cost is added into the computing cost function in the form of a weighted negative indicator function on satisfying the required delay bound, to guarantee the worst-case delay to the maximum extent.

Meanwhile, a device-related cost G_k , is generated when each task of device k is executed at the edge server for capturing the cost of task transmission, service maintenance, e.g., energy consumption, and wear and tear of the equipment [33]. The cost function is represented by

$$C_k(\mathbf{s}(t)) = \begin{cases} G_k, & \text{if } t \in \mathcal{T}_k^0 \text{ and } z_k(t) \neq 0 \\ 0, & \text{otherwise.} \end{cases} \quad (16)$$

Thus, the overall reward and cost for this network are defined as

$$R(\mathbf{s}(t)) = \sum_{k \in \mathcal{K}} R_k(\mathbf{s}(t)) \quad (17)$$

$$C(\mathbf{s}(t)) = \sum_{k \in \mathcal{K}} C_k(\mathbf{s}(t)). \quad (18)$$

Because the system evolves in a Markov way, to maximize the expectation of the long-term average reward subject to edge computing cost constraint, a CMDP problem for choosing the optimal task scheduling action in each time slot can be formulated as P0,

$$\begin{aligned} \mathbf{P0} : \max_{\pi} & \mathbb{E} \left[\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T R(\mathbf{s}(t)) \middle| \pi \right] \\ \text{s.t.} & \mathbb{E} \left[\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T C(\mathbf{s}(t)) \middle| \pi \right] \leq \bar{C} \end{aligned} \quad (19)$$

where π is a stationary policy mapping from state $\mathbf{s}(t)$ to action $\mathbf{A}(t)$. For the scenario with some urgent tasks of high priority, although the delay calculation in (11) for each task may be biased due to the reordered task computing, for one fixed policy, it can still complete the same number of tasks within the same time regardless of the serving order, i.e., the sum of the task completion profit and the sum of the task completion delay are the same. Therefore, the average performance of P0 in (19) under the above formulation is still the same as that of the scenario without reordering the tasks.

V. ALGORITHM FOR TASK SCHEDULING

If we know the environment information, i.e., the state transition probabilities $P(\mathbf{s}|\mathbf{s}', \mathbf{A}')$, the CMDP defined as P0 can be equivalently transformed to a linear programming (LP) problem [34], given by

$$\max_{\mathbf{x}} \sum_{\mathbf{s} \in \mathcal{S}} \sum_{\mathbf{A} \in \mathcal{A}(\mathbf{s})} x(\mathbf{s}, \mathbf{A}) R(\mathbf{s}) \quad (20a)$$

$$\text{s.t.} \sum_{\mathbf{s}' \in \mathcal{S}} \sum_{\mathbf{A}' \in \mathcal{A}(\mathbf{s}')} x(\mathbf{s}', \mathbf{A}') P(\mathbf{s}|\mathbf{s}', \mathbf{A}') = \sum_{\mathbf{A} \in \mathcal{A}(\mathbf{s})} x(\mathbf{s}, \mathbf{A}), \quad \forall \mathbf{s} \in \mathcal{S} \quad (20b)$$

$$\sum_{\mathbf{s} \in \mathcal{S}} \sum_{\mathbf{A} \in \mathcal{A}(\mathbf{s})} x(\mathbf{s}, \mathbf{A}) = 1 \quad (20c)$$

$$\sum_{\mathbf{s} \in \mathcal{S}} \sum_{\mathbf{A} \in \mathcal{A}(\mathbf{s})} x(\mathbf{s}, \mathbf{A}) C(\mathbf{s}) \leq \bar{C}. \quad (20d)$$

In (20), variable $x(\mathbf{s}, \mathbf{A})$ can be interpreted as the long-run visiting probability to the state-action pair (\mathbf{s}, \mathbf{A}) , $R(\mathbf{s})$ is the reward under state \mathbf{s} , and $C(\mathbf{s})$ is the generated cost under state \mathbf{s} . Constraint (20b) is the global balance equations for the steady-state probabilities, (20c) implies that the sum of all state-action pairs' probabilities is equal to 1, (20d) shows that the edge computing cost should be limited by cost budget \bar{C} . For brevity, time index t is written as the subscript in the following. Different from the deterministic policy in traditional MDP, the optimality can be achieved only when the policy is a mapping from a state to a distribution of available actions in the CMDP, defined by

$$P(\mathbf{A}_t = \mathbf{A} | \mathbf{s}_t = \mathbf{s}) = \frac{x(\mathbf{s}, \mathbf{A})}{\sum_{\mathbf{A}' \in \mathcal{A}(\mathbf{s})} x(\mathbf{s}, \mathbf{A}')}. \quad (21)$$

The LP problem in (20) can be solved in general via methods such as simplex algorithm and interior-point algorithm [34]. However, it is impractical to apply an LP-based algorithm in this

scenario since the transition probabilities (i.e., $P(\mathbf{s}'|\mathbf{s}, \mathbf{A})$) cannot be known a priori, and the high dimensional state representation leads to an exponentially increasing space complexity for $x(\mathbf{s}, \mathbf{A})$.

A. Task Scheduling without Edge Computing Cost Constraint

When the state transition probabilities are not accessible, reinforcement learning (RL) can be leveraged to gradually extract the optimal policy by interacting with the environment. However, the constrained RL is an elusive challenge due to the difficulty of transferring from value-based formulation in RL to the LP-based formulation in CMDP. In order to find a solution to the cost-constrained task scheduling problem, we first design an algorithm to deal with the case in which the average edge computing cost constraint is removed. The simplified problem for maximizing the long-term average reward is then approximated by maximizing the accumulated discounted reward [35], i.e.,

$$\mathbf{P1} : \max_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(\mathbf{s}_t) \middle| \pi \right] \quad (22)$$

where γ (< 1) is the discount factor to represent that the learning agent pays more attention to the current reward than the future reward. The larger γ , the closer the optimal solution of P1 approaches the undiscounted one.

When the state transition probabilities are unknown, Q-learning is a common model-free RL method which can learn the optimal policy via successively updating its estimation about the Q-values based on the interaction with the external environment [36]. The Q-values can be interpreted as an evaluation for how good a state-action pair is. Therefore, when the estimation about the Q-values is accurate, the action with the largest Q-value, i.e., $\arg \max_{\mathbf{A} \in \mathcal{A}(\mathbf{s})} Q(\mathbf{s}, \mathbf{A})$, should be taken when the learning agent observes state \mathbf{s} . The traditional tabular Q-learning maintains the estimation for Q-values by constructing a Q-value table for all state-action pairs with size $\sum_{\mathbf{s} \in \mathcal{S}} |\mathcal{A}(\mathbf{s})|$, where $|\cdot|$ is the set cardinality. In each decision epoch, an action is chosen probabilistically based on the current state \mathbf{s}_t and the Q-value table, and the probability for choosing action \mathbf{A}_t is expressed as

$$P(\mathbf{A}_t | \mathbf{s}_t) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(\mathbf{s}_t)|}, & \text{if } \mathbf{A}_t = \arg \max_{\mathbf{A}} Q(\mathbf{s}_t, \mathbf{A}) \\ \frac{\epsilon}{|\mathcal{A}(\mathbf{s}_t)|}, & \text{otherwise} \end{cases} \quad (23)$$

where ϵ is a tradeoff coefficient balancing between exploitation (choosing the most promising action) and exploration (trying a random action). After the action, a reward, $R(\mathbf{s}_t)$, is fed back and the next state \mathbf{s}_{t+1} is observed. Based on the full backup equation, only the Q-value for state-action pair $(\mathbf{s}_t, \mathbf{A}_t)$ in the Q-table is updated [37], given by

$$Q(\mathbf{s}_t, \mathbf{A}_t) = (1-\beta)Q(\mathbf{s}_t, \mathbf{A}_t) + \beta \left(R(\mathbf{s}_t) + \gamma \max_{\mathbf{A} \in \mathcal{A}(\mathbf{s}_{t+1})} Q(\mathbf{s}_{t+1}, \mathbf{A}) \right) \quad (24)$$

where β is the learning rate which balances the learning step length and the learning accuracy. Such an action selection and Q-value updating mechanism can guarantee that the optimal policy can be obtained with probability 1 [37]. However, the traditional tabular Q-learning is considerably space-costly for recording the Q-values of all possible state-action pairs, and has the curse of dimension problem when it is applied in our task scheduling problem with multi-dimensional state. To overcome the high space complexity, we adopt a DQL method [38] instead, to approximate Q-values via a DNN, i.e., $Q(\cdot; \theta)$, where θ is the network parameter consisting of the weights and biases for connecting the neural units of different layers in the DNN. The detailed DQL for task scheduling without the edge computing cost constraint is given in Algorithm 1.

The proposed algorithm is executed at the beginning of each slot by the LTE-U BS which is able to observe the system states via the inherent reporting scheme of LTE. There are two DQNs, i.e., the evaluation network and the target network, and their parameters are denoted by θ and θ' , respectively. The evaluation network can directly obtain an evaluation Q-value for a state-action pair via a forwarding calculation. Alternatively, the target Q-value for a state-action pair sample $(\mathbf{s}_j, \mathbf{A}_j)$ can be estimated via the full backup equation in (24), i.e.,

$$y_j = R_j + \gamma \max_{\mathbf{A} \in \mathcal{A}(\mathbf{s}_{j+1})} (Q'(\mathbf{s}_{j+1}, \mathbf{A}; \theta')) \quad (25)$$

where R_j is the reward feedback under state-action pair $(\mathbf{s}_j, \mathbf{A}_j)$, \mathbf{s}_{j+1} is the following state, and the Q-value for \mathbf{s}_{j+1} is derived from the forwarding calculation based on the target network. The action selection in DQL is the same as in (23), where the Q-value calculation is based on the evaluation network. The exploration probability ϵ gradually decays with a rate ϕ to transfer its concern from exploration to exploitation. After the action selection, a transition sample including the current state, current action, reward feedback, and the next state, is stored in the replay memory. The transition samples in the replay memory are randomly selected for the gradient-

Algorithm 1 DQL algorithm for task scheduling**Initialize:**Set an empty replay memory Ψ with capacity M ;Randomly initialize the parameter of the evaluation DQN (Q) as θ ;Initialize the parameter of target DQN (Q') as $\theta' = \theta$;Initialize $\epsilon = 1$ and state s_0 ;**while** $t < T$ **do****if** $t \in \mathcal{E}$ **then**

Choose action based on (23);

Execute the action and observe the reward R_t and the next state s_{t+1} ;**else**Set $A_t = 0$;**end if**Decay the exploration probability $\epsilon \leftarrow \max\{\epsilon\phi, \epsilon_{\min}\}$;Store transition vector (s_t, A_t, R_t, s_{t+1}) in Ψ ;Sample random minibatch of transitions (s_j, A_j, R_j, s_{j+1}) from Ψ ;Calculate the target Q value of (s_j, A_j) , i.e., y_j , based on (25);Perform a gradient decent on θ for minimizing $|y_j - Q(s_j, A_j; \theta)|^2$; $t = t + 1$;Every N steps set $\theta' = \theta$;**end while****Output:** Evaluation DQN with parameter θ

decent-based training of the evaluation DQN for shrinking the gaps between the evaluation Q-values and the target Q-values. After every N steps, the target DQN gets updated using θ , the parameter of the trained evaluation DQN. The design of two separated DQNs reduces the oscillations and avoids the divergence of the policy. By iteratively training the evaluation DQN, the agent can gradually obtain the optimal policy [39].

B. Task Scheduling with Edge Computing Cost Constraint

Developing a solution for the constrained task scheduling is mainly inspired by [40], where the constraint can be integrated into the reward function by an appropriate weight, λ , leading to a new MDP with a revised reward and no constraint. From the perspective of the DQN formulation with parameter θ , the original problem P0 with constraint can be expressed as

$$\mathbf{P2:} \max_{\theta} V_{\pi(\theta)} \quad \text{s.t.} \quad \mathbb{E}(C(t)|\pi(\theta)) \leq \bar{C} \quad (26)$$

where $\pi(\theta)$ is the policy that the action with the maximal Q value for the current state is selected based on a DQN parameterized by θ , and $V_{\pi(\theta)} = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left(\sum_{t=0}^{\infty} R(s_t) | \pi(\theta) \right)$ is the expectation of the long-term average reward under policy $\pi(\theta)$. The objective of task scheduling

is to maximize the average reward while guaranteeing the edge computing cost constraint. The Lagrange function and the dual function of prime problem P2 are expressed, respectively, as

$$L(\boldsymbol{\theta}, \lambda) = V_{\pi(\boldsymbol{\theta})} - \lambda (\mathbb{E}(C|\pi(\boldsymbol{\theta})) - \bar{C}), \quad \lambda > 0 \quad (27)$$

$$g(\lambda) = \sup_{\boldsymbol{\theta}} \{V_{\pi(\boldsymbol{\theta})} - \lambda (\mathbb{E}(C|\pi(\boldsymbol{\theta})) - \bar{C})\}. \quad (28)$$

Because the underlying problem is a CMDP with a standard LP form in (20), it is obvious that the dual problem (i.e., $\min_{\lambda} g(\lambda)$) has the same solution as the prime problem, and the dual function is reformulated as

$$g(\lambda) = \sup_{\boldsymbol{\theta}} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{t=0}^T r(\mathbf{s}_t, \mathbf{A}_t) - \lambda \sum_{t=0}^T (C(\mathbf{s}_t, \mathbf{A}_t) - \bar{C}) \middle| \pi(\boldsymbol{\theta}) \right], \lambda > 0. \quad (29)$$

Theorem 1. *Function $g(\lambda)$ is continuous even when the scheduling policy is restricted to be deterministic policy, which means that, for each state, there is only one optimal action and random policy is not allowed.*

Proof. See Appendix A. □

Therefore, the optimal expected reward in P2 is the same as the optimal value of $g(\lambda^*)$ when the optimal Lagrange multiplier λ^* is found. Accordingly, the constrained scheduling problem can be transformed as another scheduling problem without constraint, but with a revised reward function, given by

$$R'(\mathbf{s}_t) = R(\mathbf{s}_t) - \lambda^*(C(\mathbf{s}_t) - \bar{C}). \quad (30)$$

It is noted that the revised MDP and the original MDP have the same state space, action space, and transition probabilities. In order to obtain the optimal Lagrange multiplier and the optimal scheduling policy under the average edge computing cost constraint, an iteration-based algorithm is developed based on the subgradient theory, which is presented in Algorithm 2. The key idea is to use the counted average edge computing cost to update the Lagrangian multiplier, where α_i is the step size for updating of λ , and $[\cdot]^+$ is the non-negative function.

Assumption 1. *In Algorithm 2, the counted average edge computing cost, C_i , has mean of $\mu_i = \mathbb{E}[C|\pi(\boldsymbol{\theta}_i)]$ and relatively small variance σ_i with $\sigma_i \leq M\mu_i$, where M is a large positive constant.*

Assumption 2. *The DQL algorithm without constraint in Algorithm 1 can obtain the optimal policy when any revised reward function with determined λ is given.*

For Assumption 1, if we count the average edge computing cost in a long term, the variance can be small enough according to the law of large numbers; For Assumption 2, many existing works prove that the DQL algorithm is effective in addressing such scheduling problems [41].

Theorem 2. *In Algorithm 2, when $i \rightarrow \infty$, λ converges to optimal λ^* with probability 1 under Assumptions 1 and 2.*

Proof. See Appendix B. □

Algorithm 2 CDQL algorithm for task scheduling with edge edge computing cost constraint

Initialize:

Set $i = 0$, $\lambda_i = 0$, and the step size $\alpha_i = 0.5$;

while 1 do

Set the reward function as: $R_i(\mathbf{s}) = R(\mathbf{s}) - \lambda_i C(\mathbf{s})$;

Run Algorithm 1 with reward function R_i and obtain the optimal policy $\pi(\theta_i)$;

Count the average edge computing cost C_i ;

Update λ as: $\lambda_{i+1} = [\lambda_i + (C_i - \bar{C})\alpha_i]^+$;

Update the step size as $\alpha_{i+1} = \alpha_0 / (i + 1)$;

$i = i + 1$;

end while

Output: λ_i and DQN with parameter θ

VI. REALIZATION OF THE CDQL

To enhance the performance and efficiency of the proposed CDQL-based task scheduling algorithm, we make the following customized modifications for its realization:

- **Structure of the DQN:** The output dimension in the conventional DQN is usually the same as the size of action space, in order to identify the difference for adopting different available actions. However, the available action space in our problem varies with different states and the generalized action space is 2^K , which indicates that constructing a DQN where the last output layer has the same size as the action space is unrealistic in the practical implementation. Therefore, we construct the output layer as shown in Fig. 3, where the Q-value for a state-action pair is calculated based on the action value layer, action mask, and the base value node. The base value node can be considered as a basic estimation of the average Q-value of the current state, while the action value layer with size $2K$ depicts the effect of each device's action on the overall Q-value. Specifically, the $(2k - 1)$ -th or

1
2
3 $2k$ -th node in the act-value layer represents the Q-value bias when action 0 or 1 is adopted
4 for device k . The input action makes up the action mask layer by a binary coding, and
5 extracts only the act-value of the selected action to the Q-value output. The Q-value under
6 this modified structure is calculated as
7
8

$$9 \quad Q(\mathbf{s}, \mathbf{A}; \boldsymbol{\theta}) = \text{Base_value}(\mathbf{s}; \boldsymbol{\theta}) + \text{Action_value}(\mathbf{s}; \boldsymbol{\theta}) \text{Action_mask}(\mathbf{A})^T. \quad (31)$$

- 10
11
12
13 • **Preprocessing of data samples:** The traditional memory replay of DQL records all past
14 samples of the state-action pairs, rewards, and future states. However, in our problem, there
15 are many time slots in which no action is required. In other words, regardless of the adopted
16 policy, these transitions are the same, which means that feeding these transition samples in
17 the learning process cannot provide any extra information for the policy learning but leads
18 to some redundant processing time. To utilize the data samples in a more efficient way, the
19 samples for training are compressed with form $(\mathbf{s}_t, \mathbf{A}_t, T, R_{t \sim T}, \mathbf{s}_{t+T})$, where \mathbf{s}_t is a state
20 under a decision epoch, \mathbf{A}_t is its corresponding action, T is the number of the slots between
21 the former decision epoch and the later decision epoch, \mathbf{s}_{t+T} is the state under the later
22 decision epoch, and $R_{t \sim T} = \sum_{i=t}^{t+T} \gamma^{i-t} R_i$ is the accumulated discounted rewards between slot
23 t and $t + T$. With the modified memory replay design, the target Q value for state-action
24 pair $(\mathbf{s}_t, \mathbf{A}_t)$ is calculated as
25
26
27
28
29
30
31
32
33

$$34 \quad y = R_{t \sim T} + \gamma^T \max_{\mathbf{A} \in \mathcal{A}(\mathbf{s}_{t+T})} Q(\mathbf{s}_{t+T}, \mathbf{A}; \boldsymbol{\theta}). \quad (32)$$

- 35
36
37
38 • **Pretraining for the Q-value:** The objective of training the network is to 1) approximate
39 the Q values for different state-action pairs with high accuracy and 2) learn the optimal
40 policy based on the trained DQN. However, the adopted policy is not stationary in the
41 initial stage of the DQN training and affects both input (i.e., input action) and output (i.e.,
42 output Q values) of the DQN. Such a recurrent effect of the policy makes it necessary to
43 choose a very small learning rate (e.g., 10^{-6}) to guarantee the convergence and stability of
44 the concurrent learning of both Q-value and policy, which leads to very time-consuming
45 training process. Therefore, we treat the policy as a random policy in the initial training stage
46 to reduce the impact of the policy dynamics, and set learning rate at a relatively large value
47 (e.g., 10^{-3}). At the same time, the accumulated discounted reward, $\hat{Q}(t) = \sum_{i=t-1000}^t R_i$ is
48 counted as an approximation of Q-value under the random policy. After the output Q-value
49
50
51
52
53
54
55
56
57
58
59
60

roughly approximates the real discounted accumulated reward, e.g., $0.9 < \frac{Q(s_t, \mathbf{A}_t; \theta)}{\hat{Q}(t)} < 1.1$, we reduce the learning rate to a small value to further optimize the policy.

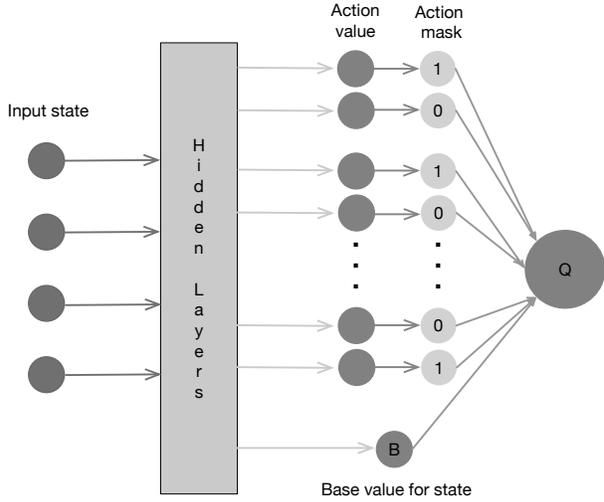


Fig. 3: The modified structure of DQN.

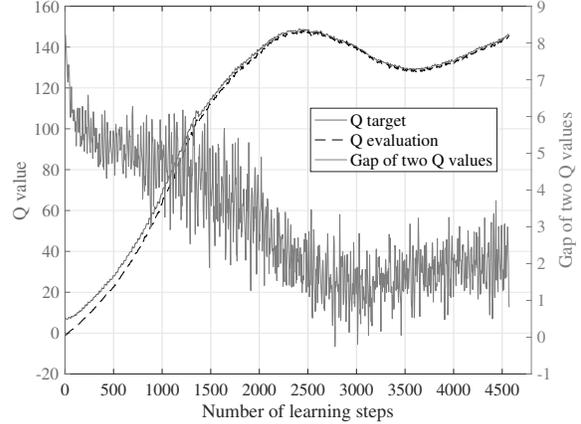


Fig. 4: Q values in the initial training stage without pretraining.

VII. PERFORMANCE EVALUATION

In this section, simulation results are presented to investigate the performance of the proposed computing offloading scheme. A python-based simulator including the LBT, data transmission, and computing is adopted for the simulation. In the simulation, we consider an LTE-U-enabled IoT network consisting of a BS with the coverage radius of 100m. There are in total 16 randomly distributed IoT devices associated with the BS. The LTE-U network is operated on the band of 5GHz and the distance-dependant path loss model is given by [42]

$$PL(R) = 38.46 \log_{10}(R) + 20 \log_{10}(R) + 0.7R \quad (33)$$

where R is the distance from the device to the LTE-U BS. Small-scale Rayleigh fading is also considered, where channel power gain follows an exponential distribution with unit mean. The noise power spectrum density is -174dBm/Hz. The dedicated bandwidth for the computing offloading service is 1.4 MHz and is equally allocated to all devices. There are 6 Wi-Fi devices with full-buffer traffic model randomly located in the network coverage area, employing the IEEE 802.11ac distributed coordination function (DCF) protocol for channel access with the same configuration in [26]. The task input data size of different devices (L_k) follows a uniform distribution from 15 Kbits to 20 Kbits. The workloads for all devices are set as 200 cycles/bit.

1
2
3 The computing capacity of the edge server and the IoT devices are 5 G cycles/s and 500 M
4 cycles/s. If not specified, the defaulted task arrival rate at each device is set as 100/s. There
5 are two backoff stages for LTE-U's LBT-based channel access and the values of the backoff
6 windows are set as 8 and 16 as the channel access priority 2 in [43]. The time duration of each
7 successful transmission phase is 10ms. The cost for computing a task at the edge server is set
8 as 5 and the edge computing cost budget is 3/ms. The duration of each slot is 1ms. The reward
9 of completing a task is set as a random value between 6 and 9 for different devices. The weight
10 for delay, η , is set as 0.1, and the discount factor is set as $\gamma = 0.99$.

11
12 The DQN is implemented based on Tensorflow [44], which is an open-source library. Unless
13 otherwise specified, the structure and setting for the DQNs are set as follows: There are three
14 hidden layers with (256, 128, 64) neural units in the DQN structure of Fig. 3. ReLu is the
15 activation function for all the hidden layers and there is no nonlinear activation function for other
16 layers. The replay memory can store $M = 5000$ transition samples and the training batch size is
17 set as 64. The parameter of the target DQN is replaced after every $N = 50$ slots. The learning
18 rates for the pretraining and training are set as 10^{-3} and 10^{-6} , respectively. Every learning
19 episode consists of 5×10^5 time slots. In one episode learning, the naive Q-learning without
20 the customized modifications is selected as the benchmark. In addition, a Lyapunov optimization
21 (LO)-based scheduling algorithm is set as the task scheduling benchmark. Because the delay
22 term in the reward function (15) depends on the stochastic channel states in the following time
23 slots, and cannot be determined in the current slot in which the decision is made. Therefore, for
24 the benchmark comparison, the delay-related term in (15) is estimated based on the environment
25 statistics. For guaranteeing the average computing cost within the computing budget, a virtual
26 queue method is applied [45].

27
28 In Fig. 4, we compare the evaluation Q-values with the target Q-values in the early training
29 phase when there is no pretraining for Algorithm 1, where both Q-values are averaged by a
30 moving window of size 10. It is shown that the Q-value of the target network is always higher
31 than that of the evaluation network in the initial training stage. Actually, such gap is mainly due
32 to their underestimations of the Q-values with the small weight initialization for the two DQNs.
33 Recall that the target Q-value is calculated by $y_j = R_j + \gamma \max_{A' \in \mathcal{A}(s_{j+1})} (Q'(s_{j+1}, A'; \theta'))$. Although
34 the future part of the target Q-value is also underestimated as that of the evaluation Q-value, the
35 target Q-value has the instant reward (usually being positive) feedback as a correction, making
36 the Q value of the target network larger than that of the evaluation network. When a small
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

learning rate is applied for the stable convergence of the DQN, many learning steps would be wasted to reduce such a gap, which makes the training inefficient. Therefore, it is essential to design an adaptive learning rate scheme for the DQN to properly estimate the Q-values in the early training stage.

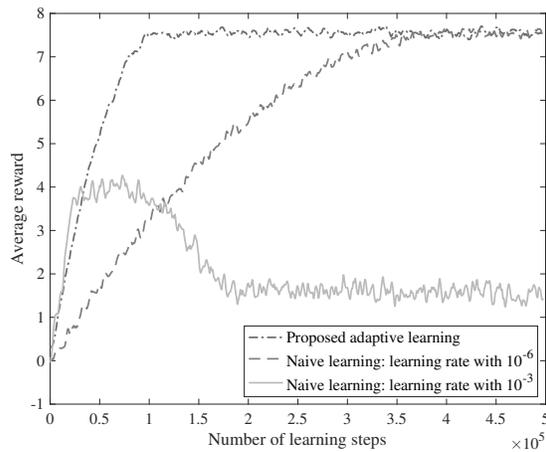


Fig. 5: Training process with different learning rate settings.

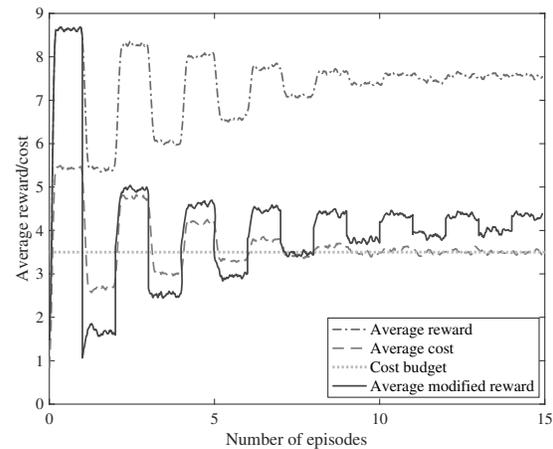


Fig. 6: Convergence of Algorithm 2.

In Fig. 5, the reward convergence process in the first episode (with $\lambda = 0$) is given under our proposed Q-learning algorithm with customized modifications and the naive Q-learning algorithm. It can be found that when the learning rate is high (e.g., 10^{-3}), the naive learning process is highly unstable and even becomes worse after training. When the learning rate is small (e.g., 10^{-6}), although the average reward can converge to a steady and high level, the learning process is considerably slow. However, in our proposed learning with customized modifications, the learning rate is adaptive, i.e., a high learning rate in the pretraining and a low learning rate after the pretraining, the achieved average performance is much better and stable, and takes only less than half time of the learning with fix rate of 10^{-6} , validating the effectiveness of our proposed modifications.

The convergence process of Algorithm 2 is shown in Fig. 6, where the average modified reward, average reward, and average cost are given. When the agent keeps receiving the cost feedback via the interaction with the environment and adapting the weight for edge computing cost, the average cost gradually matches the cost budget, and all average modified reward, average reward and average cost become stable, validating the effectiveness of Algorithm 2 for guaranteeing the cost budget. Note that pretraining is applied only in the first episode, while the learning in the subsequent episodes is directly trained based on the DNN parameters in the former episode.

However, we can see the learning process can converge quickly when the target (accumulative revised reward function) is changed. It is because that some learned features in the hidden layers can be shared as common knowledge among the learning models with different objectives. Such a phenomena can be interpreted as the advantage of transfer learning [46].

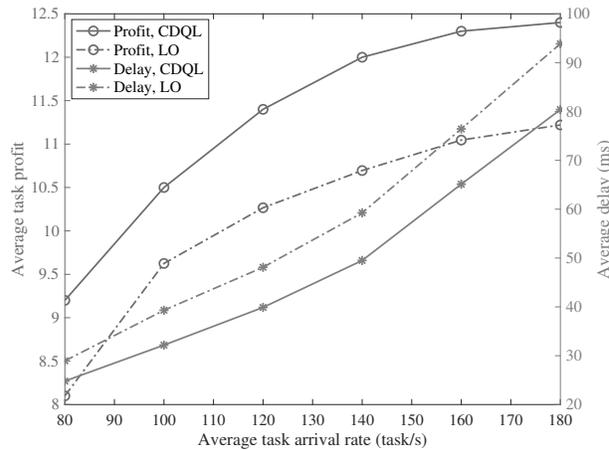


Fig. 7: Performance comparison between CDQL and LO.

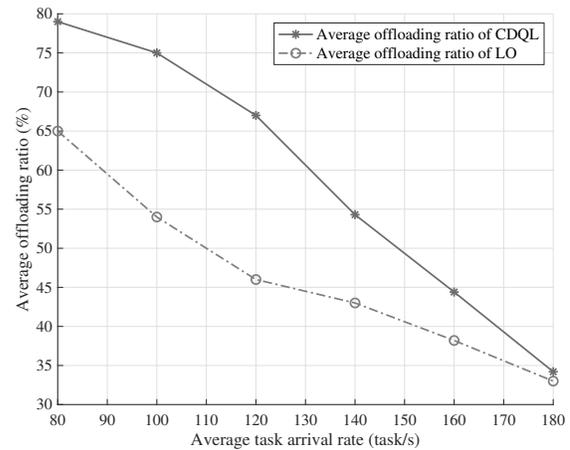


Fig. 8: Average offloading ratio versus task arrival rates.

In Fig. 7, we compare the performance of our proposed algorithm with that of the LO algorithm. One can notice that, for both algorithms with the same average edge computing cost constraint, as the task arrival rate increases, the average task profit increases with a decreasing marginal gain due to the gradually saturated communication and computing capability of the system, while the average delay increases considerably. However, our proposed CDQL algorithm always outperforms the LO algorithm in terms of both average profit and delay. The reasons are two-folded: Firstly, the LO algorithm suffers from the inaccuracy of the reward based on the estimated task completion delay; Secondly, the main idea of the LO algorithm is to utilize the relationship between the queue drift and the adopted actions, thus cannot exploit the environment dynamics, while our proposed algorithm can learn the channel state evolution and the task arrival pattern, leading to higher task completion profit and lower delay.

The average offloading ratios for both CDQL algorithm and LO algorithm with different task arrival rates are shown in Fig. 8, where both algorithms maintain the same edge computing cost budget. Due to more intelligent scheduling for task offloading, the proposed CDQL algorithm achieves a higher offloading ratio, leading to more efficient utilization of the edge computing resources, especially when the task arrival rate is small. However, when the task arrival rate increases, the utilization of edge computing resources becomes gradually saturated.

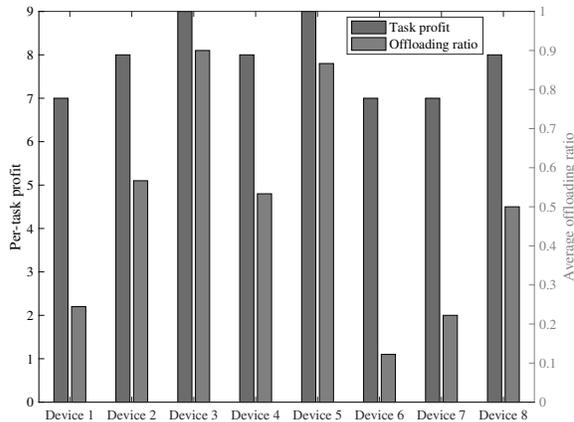


Fig. 9: Offloading ratio and task profit for different IoT devices.

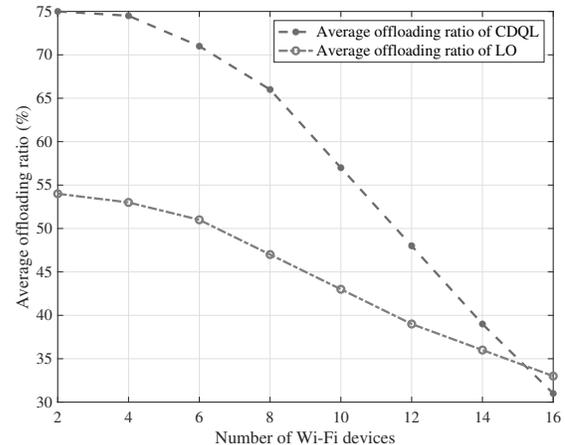


Fig. 10: Average offloading ratio versus number of Wi-Fi devices.

Figure 9 compares the average offloading ratio with per-task completion profit for randomly chosen 8 IoT devices in the LTE-U network when the task arrival rate is set as 160/s. It is found out that a device with a higher task profit can have a larger offloading ratio. Because the computing capability of the edge server is much powerful than the IoT device, processing task at the edge server is usually a more efficient way than local computing. By using our proposed task scheduling algorithm, the system can intelligently allocate more computing resources (i.e., higher offloading ratio) to those devices with higher profit because the proposed algorithm can exploit the diversity among different devices, the system state and the channel availability information.

Figure 10 illustrates the average offloading ratios when the LTE-U-enabled IoT coexisting with different number of Wi-Fi devices. It can be found out when the number of Wi-Fi devices is small, the offloading ratio maintains a stable level even if the number of Wi-Fi devices increases. It is because under such a condition the edge computing cost budget is the main bottleneck for offloading tasks. When the number of Wi-Fi devices further increases, the unlicensed channel becomes more crowded and offloading tasks to the edge server may lead to a high transmission delay. Therefore, offloading may not be a more efficient selection and the offloading ratio decreases. However, it can be seen that the LO algorithm is not as intelligent as our proposed algorithm in terms of adapting the offloading ratio to the network environment.

VIII. CONCLUSION

In this paper, we propose an RL-based method to solve the task scheduling problem for the computing offloading in an LTE-U-enabled IoT network. We aim to maximize the average reward that integrates both task completion profit and task delay, while guaranteeing average

edge computing cost constraint. We formulate the scheduling problem as a CMDP to capture the dynamics of the unlicensed channels' availability and the task traffic. To deal with the exponentially increased space complexity due to the high state dimensions in conventional Q-learning, we utilize a DQL-based method to approximate Q values of different state-action pairs, instead of using the traditional tabular method. Based on the duality theory, a CDQL framework is proposed to integrate the average edge computing cost constraint into the DQL framework. Specifically, we define a revised reward function that combines the original reward and the edge computing cost weighted by an appropriately determined Lagrange multiplier, such that the constraint can be removed in the CDQL framework. In addition, a sub-gradient-based method is proposed to search for the optimal Lagrange multiplier, and the convergence of the searching algorithm is proved. Simulation results demonstrate that our proposed algorithm considerably improves the system performance in terms of the task completion profit and the task completion delay. For the future work, we will extend our work to the scenario where each device can learn the optimal task scheduling policy in a distributed way via multi-agent reinforcement learning.

APPENDIX A

To prove that $g(\lambda)$ is a continuous function, we should prove that for any λ and any ϵ , there exists δ such that, as long as $\|\lambda' - \lambda\| \leq \delta$, we have $|g(\lambda) - g(\lambda')| \leq \epsilon$, where $\|\cdot\|$ and $|\cdot|$ are Euclidean norm function and absolute-value function, respectively. The available policy set is denoted as $\mathcal{L} = \{1, 2, \dots, L\}$, where $L \leq |\mathcal{A}|^{|\mathcal{S}|}$ is the number of all available deterministic policies, and is a very large but finite positive integer. The optimal policy for λ is denoted as l_λ , i.e.,

$$l_\lambda = \arg \max_{l \in \mathcal{L}} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[\sum_{t=0}^T R(\mathbf{s}_t) - \lambda \sum_{t=0}^T (C(\mathbf{s}_t) - \bar{C}) \middle| \pi(\boldsymbol{\theta}) \right]. \quad (34)$$

For each deterministic policy $l \in \mathcal{L}$, function $L(\boldsymbol{\theta}(l), \lambda)$ in (27) with a given $\boldsymbol{\theta}$ is continuous with λ , where $\boldsymbol{\theta}(l)$ is the parameter for policy l . Therefore, we conclude that there exists δ_l such that, for any ϵ , as long as $\|\lambda' - \lambda\| \leq \delta_l$, $|L(\boldsymbol{\theta}(l), \lambda) - L(\boldsymbol{\theta}(l), \lambda')| \leq \epsilon$. It is natural to postulate that we choose $\delta = \min_{l \in \mathcal{L}} \delta_l$. For any λ' that satisfies $\|\lambda' - \lambda\| \leq \delta$, we can have

$$\begin{aligned} g(\lambda) - g(\lambda') &= L(\boldsymbol{\theta}(l_\lambda), \lambda) - L(\boldsymbol{\theta}(l_{\lambda'}), \lambda') \stackrel{(a)}{\leq} L(\boldsymbol{\theta}(l_\lambda), \lambda') + \epsilon - L(\boldsymbol{\theta}(l_{\lambda'}), \lambda') \\ &\stackrel{(b)}{\leq} L(\boldsymbol{\theta}(l_\lambda), \lambda') + \epsilon - L(\boldsymbol{\theta}(l_\lambda), \lambda') \leq \epsilon. \end{aligned} \quad (35)$$

In (35), inequality (a) holds because $L(\boldsymbol{\theta}(l_\lambda), \lambda) \leq L(\boldsymbol{\theta}(l_\lambda), \lambda') + \epsilon$ as $\|\lambda' - \lambda\| \leq \delta \leq \delta_{l_\lambda}$ and inequality (b) holds because $\boldsymbol{\theta}(l_{\lambda'})$ is the optimal policy parameter over any other parameter

under Lagrangian multiplier λ' .

Further,

$$\begin{aligned} g(\lambda) - g(\lambda') &= L(\boldsymbol{\theta}(l_\lambda), \lambda) - L(\boldsymbol{\theta}(l_{\lambda'}), \lambda') \stackrel{(c)}{\geq} L(\boldsymbol{\theta}(l_\lambda), \lambda) - (L(\boldsymbol{\theta}(l_{\lambda'}), \lambda) + \epsilon) \\ &\stackrel{(d)}{\geq} L(\boldsymbol{\theta}(l'_\lambda), \lambda) - (L(\boldsymbol{\theta}(l_{\lambda'}), \lambda) + \epsilon) \geq -\epsilon \end{aligned} \quad (36)$$

where inequality (c) holds because $L(\boldsymbol{\theta}(l_{\lambda'}), \lambda') \leq L(\boldsymbol{\theta}(l_{\lambda'}), \lambda) + \epsilon$ as $\|\lambda' - \lambda\| \leq \delta \leq \delta_{l_{\lambda'}}$ and inequality (d) holds because $\boldsymbol{\theta}(l_\lambda)$ is the optimal policy parameter over any other parameter under Lagrangian multiplier λ .

As a result, we can conclude that, as long as $\|\lambda' - \lambda\| \leq \delta = \min_{l \in \mathcal{L}} \delta_l$, $|g(\lambda) - g(\lambda')| \leq \epsilon$, i.e., $g(\lambda)$ is a continuous function.

APPENDIX B

Lemma 1. Suppose $\{X_n, n \geq 1\}$ are independent non-negative random variables with $E(X_n) = \mu_n$, $Var(X_n) = \sigma_n^2$. Define for $n \geq 1$, $S_n = \sum_{i=1}^n X_i$, and suppose that $\sum \mu_i = \infty$ and $\sigma_n^2 \leq M\mu_n$ for some $M > 0$ and all n . Then sequence $S_n/\mathbb{E}(S_n)$ converges to 1 in probability.

Proof. According to Chebychev's inequality, for any $\epsilon > 0$, we have

$$P\{|S_n/\mathbb{E}(S_n) - 1| \geq \epsilon\} \leq \frac{Var(S_n/\mathbb{E}(S_n))}{\epsilon^2} = \frac{Var(S_n)}{\mathbb{E}^2(S_n)\epsilon^2}. \quad (37)$$

For $S_n = \sum_{i=1}^n X_i$, we have

$$P\{|S_n/\mathbb{E}(S_n) - 1| \geq \epsilon\} \leq \frac{Var\left(\sum_{i=1}^n X_i\right)}{\mathbb{E}^2\left(\sum_{i=1}^n X_i\right)\epsilon^2} = \frac{\sum_i \sigma_i^2}{\left(\sum_i \mu_i\right)^2 \epsilon^2} \leq \frac{M \sum_i \mu_i}{\left(\sum_i \mu_i\right)^2 \epsilon^2} = \frac{M}{\sum_i \mu_i \epsilon^2}. \quad (38)$$

Therefore, when $n \rightarrow +\infty$, we can obtain that

$$\lim_{n \rightarrow +\infty} P\{|S_n/\mathbb{E}(S_n) - 1| \geq \epsilon\} = 0. \quad (39)$$

□

Definition 1. A subgradient of convex function f at \mathbf{x} is any vector \mathbf{z} that satisfies the inequality $f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{z}^T(\mathbf{y} - \mathbf{x})$ for all $\mathbf{y} \in \text{dom} f$, where $\text{dom} f$ is the domain of function f [47].

Lemma 2. The subgradient of function $g(\lambda)$ can be denoted as $h(\boldsymbol{\theta}_\lambda) = \mathbb{E}[C|\pi(\boldsymbol{\theta}_\lambda)] - \bar{C}$, where $\boldsymbol{\theta}_\lambda$ is the value of $\boldsymbol{\theta}$ when $L(\boldsymbol{\theta}, \lambda)$ gets the supremum given λ , i.e., $\boldsymbol{\theta}_\lambda = \arg \max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \lambda)$.

Proof. For any $\lambda' > \mathbf{0}$, we have

$$\begin{aligned}
g(\lambda') &= \max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}, \lambda') = \max_{\boldsymbol{\theta}} [V_{\pi(\boldsymbol{\theta})} + \lambda' h(\boldsymbol{\theta})] \\
&\geq V_{\pi(\boldsymbol{\theta}_\lambda)} + \lambda' h(\boldsymbol{\theta}_\lambda) = V_{\pi(\boldsymbol{\theta}_\lambda)} + \lambda h(\boldsymbol{\theta}_\lambda) + (\lambda' - \lambda) h(\boldsymbol{\theta}_\lambda) \\
&= g(\lambda) + h(\boldsymbol{\theta}_\lambda)(\lambda' - \lambda). \tag{40}
\end{aligned}$$

□

Assume that the subgradient of function $g(\lambda)$ is bounded, i.e., $\|h(\lambda)\|_2 \leq G$. We prove the convergence by showing that the distance between optimal λ^* for the minimal $g(\lambda)$ and λ_i converges to zero. In the i -th iteration of Algorithm 2, the counted subgradient function is defined as $\bar{h}_i = C_i - \bar{C}$, and the measurement error for the expectation of average cost is defined as $n_i = C_i - \mathbb{E}[C|\pi(\boldsymbol{\theta}_i)]$. Thus the true subgradient in the i -th iteration of Algorithm 2 is denoted as $h_i = \mathbb{E}[C|\pi(\boldsymbol{\theta}_i)] - \bar{C} = \bar{h}_i - n_i$. Accordingly, we have

$$\begin{aligned}
\|\lambda_{i+1} - \lambda^*\|_2^2 &= \left\| [\lambda_i - \alpha_i \bar{h}_i]^+ - \lambda^* \right\|_2^2 \leq \|\lambda_i - \alpha_i \bar{h}_i - \lambda^*\|_2^2 = \|\lambda_i - \alpha_i (h_i + n_i) - \lambda^*\|_2^2 \\
&\leq \|\lambda_i - \lambda^*\|_2^2 - 2\alpha_i h_i (\lambda_i - \lambda^*) - 2\alpha_i n_i (\lambda_i - \lambda^*) + \alpha_i^2 \|h_i + n_i\|_2^2 \\
&\leq \|\lambda_i - \lambda^*\|_2^2 - 2\alpha_i (g(\lambda_i) - g^*) - 2\alpha_i n_i (\lambda_i - \lambda^*) + \alpha_i^2 \|h_i + n_i\|_2^2. \tag{41}
\end{aligned}$$

Applying the inequality above recursively, we have

$$\|\lambda_{i+1} - \lambda^*\|_2^2 \leq \|\lambda_1 - \lambda^*\|_2^2 - 2 \sum_{k=1}^i \alpha_k (g(\lambda_k) - g^*) + \sum_{k=1}^i \alpha_k^2 \|h_k + n_k\|_2^2 - \sum_{k=1}^i 2\alpha_k n_k (\lambda_k - \lambda^*). \tag{42}$$

Based on the fact that $\|\lambda_{i+1} - \lambda^*\|_2^2 \geq 0$ and $\|\lambda_1 - \lambda^*\|_2$ must be bounded by some large number which is denoted by M , we have

$$2 \sum_{k=1}^i \alpha_k (g(\lambda_k) - g^*) \leq M^2 + \sum_{k=1}^i \alpha_k^2 \|h_k + n_k\|_2^2 + z_i \tag{43}$$

where $z_i = \sum_{k=1}^i 2\alpha_k n_k (\lambda_k - \lambda^*)$. Then denoting the best result of $g(\lambda)$ before the $i - 1$ iterations as g_i^{best} , i.e., $g_i^{\text{best}} = \max_{k < i} g(\lambda_k)$, we have

$$\sum_{k=1}^i \alpha_k (g(\lambda_k) - g^*) \geq \left(\sum_{k=1}^i \alpha_k \right) \min_{k < i} (g(\lambda_k) - g^*) = \left(\sum_{k=1}^i \alpha_k \right) (g_i^{\text{best}} - g^*). \tag{44}$$

We can further obtain the inequality

$$g_i^{\text{best}} - g^* \leq \frac{M^2 + \sum_{k=1}^i \alpha_k^2 \|h_k + n_k\|_2^2 + z_i}{\sum_{k=1}^i \alpha_k} = \frac{M^2 + \sum_{k=1}^i \alpha_k^2 \|h_k + n_k\|_2^2}{\sum_{k=1}^i \alpha_k} + \frac{\sum_{k=1}^i 2\alpha_k (h_k - \bar{h}_k) (\lambda_k - \lambda^*)}{\sum_{k=1}^i \alpha_k}. \tag{45}$$

When we choose α_i as square summable but not summable, i.e., $\sum_i \alpha_i^2 < \infty$, and $\sum_i \alpha_i = \infty$,

it is obvious that the first term $\frac{M^2 + \sum_{k=1}^i \alpha_k^2 \|h_i + n_i\|_2^2}{\sum_{k=1}^i \alpha_k}$ converges to 0 in probability. The second term

is denoted as $\frac{\sum_{k=1}^i 2\alpha_k (h_k - \bar{h}_k)(\lambda_k - \lambda^*)}{\sum_{k=1}^i \alpha_k} = \frac{\sum_{k=1}^i 2\alpha_k h_k (\lambda_k - \lambda^*)}{\sum_{k=1}^i \alpha_k} - \frac{\sum_{k=1}^i 2\alpha_k \bar{h}_k (\lambda_k - \lambda^*)}{\sum_{k=1}^i \alpha_k} = \sum_k a_k - \sum_k b_k$, where $\sum_k a_k = \mathbb{E} \left(\sum_k b_k \right)$. Therefore, based on Lemma 2, the term in (45) converges to 0 in probability.

REFERENCES

- [1] H. He, H. Shan, A. Huang, Q. Ye, and W. Zhuang, "Reinforcement learning-based computing and transmission scheduling for LTE-U-enabled IoT," in *Proc. IEEE GLOBECOM*, Dec. 2018, pp. 1–6.
- [2] Q. Ye and W. Zhuang, "Distributed and adaptive medium access control for Internet-of-Things-enabled mobile networks," *IEEE Internet Things J.*, vol. 4, no. 2, pp. 446–460, Apr. 2017.
- [3] K. Zhang, J. Ni, K. Yang, X. Liang, J. Ren, and X. Shen, "Security and privacy in smart city applications: Challenges and solutions," *IEEE Commun. Mag.*, vol. 55, no. 1, pp. 122–129, Jan. 2017.
- [4] D. Puthal, S. P. Mohanty, V. P. Yanambaka, and E. Kougianos, "Poah: A novel consensus algorithm for fast scalable private blockchain for large-scale IoT frameworks," *arXiv preprint arXiv:2001.07297*, Jan. 2020.
- [5] Z. Qin, F. Y. Li, G. Y. Li, J. A. McCann, and Q. Ni, "Low-power wide-area networks for sustainable IoT," *IEEE Wireless Commun.*, vol. 26, no. 3, pp. 140–145, Jan. 2019.
- [6] R. S. Sinha, Y. Wei, and S.-H. Hwang, "A survey on LPWA technology: LoRa and NB-IoT," *ICT Express*, vol. 3, no. 1, pp. 14–21, Mar. 2017.
- [7] MulteFire Alliance, "Multefire release 1.1 technical overview white paper," Dec. 2018.
- [8] S. Tatesh, *eLTE brings Yangshan port into automation era*, Huawei Technologies Co., Ltd., 2018.
- [9] C. Long, Y. Cao, T. Jiang, and Q. Zhang, "Edge computing framework for cooperative video processing in multimedia IoT systems," *IEEE Trans. Multimedia*, vol. 20, no. 5, pp. 1126–1139, May 2018.
- [10] J. Ren, H. Guo, C. Xu, and Y. Zhang, "Serving at the edge: A scalable IoT architecture based on transparent computing," *IEEE Network*, vol. 31, no. 5, pp. 96–105, Aug. 2017.
- [11] X. Chen, Q. Shi, L. Yang, and J. Xu, "Thriftyedge: Resource-efficient edge computing for intelligent IoT applications," *IEEE Network*, vol. 32, no. 1, pp. 61–65, Jan. 2018.
- [12] G. Hattab and D. Cabric, "Distributed wideband sensing-based architecture for unlicensed massive IoT communications," *arXiv preprint arXiv:1812.06238*, Dec. 2018.
- [13] J. Ho and M. Jo, "Offloading wireless energy harvesting for IoT devices on unlicensed bands," *IEEE Internet Things J.*, vol. 6, no. 2, pp. 3663–3675, Apr. 2019.
- [14] H. Zhang, B. Di, K. Bian, and L. Song, "IoT-U: Cellular Internet-of-things networks over unlicensed spectrum," *IEEE Trans. Wireless Commun.*, vol. 18, no. 5, pp. 2477–2492, May 2019.

- 1
2
3 [15] S. Baidya and M. Levorato, "Edge-assisted content and computation-driven dynamic network selection for
4 real-time services in the urban IoT," in *2017 INFOCOM WKSHPs*, May 2017, pp. 796–801.
- 5 [16] Y. Chen, N. Zhang, Y. Zhang, and X. Chen, "Dynamic computation offloading in edge computing for Internet
6 of Things," *IEEE Internet Things J.*, vol. 6, no. 3, pp. 4242–4251, Jun. 2019.
- 7 [17] Y. Chen, N. Zhang, Y. Zhang, X. Chen, W. Wu, and X. S. Shen, "Energy efficient dynamic offloading in
8 mobile edge computing for Internet of things," *IEEE Trans. Cloud Comput.*, pp. 1–1, Feb. 2019.
- 9 [18] X. Lyu, W. Ni, H. Tian, R. P. Liu, X. Wang, G. B. Giannakis, and A. Paulraj, "Optimal schedule of mobile
10 edge computing for Internet of Things using partial information," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 11,
11 pp. 2606–2615, Nov. 2017.
- 12 [19] N. Cheng, F. Lyu, W. Quan, C. Zhou, H. He, W. Shi, and X. Shen, "Space/aerial-assisted computing offloading
13 for IoT applications: A learning-based approach," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 5, pp. 1117–1129,
14 May 2019.
- 15 [20] Y. Liu, H. Yu, S. Xie, and Y. Zhang, "Deep reinforcement learning for offloading and resource allocation in
16 vehicle edge computing and networks," *IEEE Trans. Veh. Technol.*, vol. 68, no. 11, pp. 11 158–11 168, Nov.
17 2019.
- 18 [21] T. Q. Dinh, Q. D. La, T. Q. S. Quek, and H. Shin, "Learning for computation offloading in mobile edge
19 computing," *IEEE Trans. Commun.*, vol. 66, no. 12, pp. 6353–6367, Dec. 2018.
- 20 [22] H. Habibzadeh, Z. Qin, T. Soyata, and B. Kantarci, "Large-scale distributed dedicated-and non-dedicated smart
21 city sensing systems," *IEEE Sensors J.*, vol. 17, no. 23, pp. 7649–7658, Dec. 2017.
- 22 [23] A. M. Haubenwaller and K. Vandikas, "Computations on the edge in the Internet of Things," *Procedia
23 Computer Science*, vol. 52, pp. 29 – 34, May 2015.
- 24 [24] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The
25 communication perspective," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 4, pp. 2322–2358, Aug. 2017.
- 26 [25] L. Yang, J. Cao, Y. Yuan, T. Li, A. Han, and A. Chan, "A framework for partitioning and execution of data
27 stream applications in mobile cloud computing," *ACM SIGMETRICS Performance Evaluation Review*, vol. 40,
28 no. 4, pp. 23–32, Mar. 2013.
- 29 [26] H. He, H. Shan, A. Huang, L. X. Cai, and T. Q. S. Quek, "Proportional fairness-based resource allocation for
30 LTE-U coexisting with Wi-Fi," *IEEE Access*, vol. 5, pp. 4720–4731, Apr. 2017.
- 31 [27] 3GPP TR 36.889, "Study on licensed-assisted access to unlicensed spectrum," May 2015.
- 32 [28] X. Chen, "Decentralized computation offloading game for mobile cloud computing," *IEEE Trans. Parallel
33 Distrib. Syst.*, vol. 26, no. 4, pp. 974–983, Apr. 2015.
- 34 [29] M. Hirzallah, M. Krunz, and Y. Xiao, "Harmonious cross-technology coexistence with heterogeneous traffic
35 in unlicensed bands: Analysis and approximations," *IEEE Trans. Cognitive Commun. and Net.*, vol. 5, no. 3,
36 pp. 690–701, Aug. 2019.
- 37 [30] G. J. Sutton, R. P. Liu, and Y. J. Guo, "Delay and reliability of load-based listen-before-talk in LAA," *IEEE
38 Access*, vol. 6, pp. 6171–6182, Dec. 2017.
- 39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

- 1
2
3 [31] M. Hirzallah, W. Afifi, and M. Krunz, "Full-duplex-based rate/mode adaptation strategies for Wi-Fi/LTE-U
4 coexistence: A POMDP approach," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 1, pp. 20–29, Nov. 2016.
- 5 [32] Y. Mao, J. Zhang, and K. B. Letaief, "Dynamic computation offloading for mobile-edge computing with
6 energy harvesting devices," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 12, pp. 3590–3605, Dec. 2016.
- 7 [33] Y. Wang, M. Sheng, X. Wang, L. Wang, and J. Li, "Mobile-edge computing: Partial computation offloading
8 using dynamic voltage scaling," *IEEE Trans. Commun.*, vol. 64, no. 10, pp. 4268–4282, Oct. 2016.
- 9 [34] H. C. Tijms, *A First Course in Stochastic Models*. John Wiley and Sons, 2003.
- 10 [35] J. N. Tsitsiklis and B. Van Roy, "On average versus discounted reward temporal-difference learning," *Machine
11 Learning*, vol. 49, no. 2-3, pp. 179–191, Nov. 2002.
- 12 [36] C. J. Watkins and P. Dayan, "Q-learning," *Machine Learning*, vol. 8, no. 3-4, pp. 279–292, May 1992.
- 13 [37] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An introduction*. Cambridge, MA: MIT Press, 2011.
- 14 [38] V. Mnih *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, p.
15 529, Feb. 2015.
- 16 [39] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double Q-learning," in *Proc. of
17 Conf. Artificial Intelligence (AAAI)*, Mar. 2016.
- 18 [40] A. M. Makowski and A. Schwartz, "Optimal index policies for MDPs with a constraint," in *Proc. of IEEE
19 CDC*, 1991, pp. 471–476.
- 20 [41] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari
21 with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, Dec. 2013.
- 22 [42] F. Liu, E. Bala, E. Erkip, M. C. Beluri, and R. Yang, "Small-cell traffic balancing over licensed and unlicensed
23 bands," *IEEE Trans. Veh. Technol.*, vol. 64, no. 12, pp. 5850–5865, Dec. 2015.
- 24 [43] 3GPP TS 36.213, "Evolved universal terrestrial radio access (E-UTRA); physical layer procedures (Release
25 14)," 2019.
- 26 [44] M. Abadi *et al.*, "Tensorflow: A system for large-scale machine learning," in *Proc. OSDI*, Nov. 2016, pp.
27 265–283.
- 28 [45] N. Lu, B. Ji, and B. Li, "Age-based scheduling: Improving data freshness for wireless real-time traffic," in
29 *Proceedings of the ACM International Symposium on Mobile Ad Hoc Networking and Computing (MobiHoc)*,
30 2018, pp. 191–200.
- 31 [46] Y. Bengio, "Deep learning of representations for unsupervised and transfer learning," in *Proceedings of ICML
32 Workshop on Unsupervised and Transfer Learning*, Jul. 2011, pp. 17–37.
- 33 [47] S. Boyd and J. Park, "Subgradient methods," *Lecture notes in EE364b, Stanford University*, 2013.
- 34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60