# Nonlinear Online Incentive Mechanism Design in Edge Computing Systems with Energy Budget

Gang Li, *Student Member, IEEE*, Jun Cai, *Senior Member, IEEE*, Xianfu Chen, and Zhou Su, *Senior Member, IEEE*

**Abstract**—In this paper, we consider task offloading in edge computing systems, where tasks are offloaded by the base station to resourceful mobile users. With the consideration of unique characteristics in practical edge computing systems, such as dynamic arrival of computation tasks, and energy constraints at battery-powered mobile users, we formulate an incentive mechanism design problem by jointly optimizing task offloading decisions, and allocation of both communications (i.e., power and bandwidth), and computation resources. In order to tackle the nonlinear issue in the designed mechanism, a novel online incentive mechanism is proposed. We first convert the original mechanism design problem into several one-shot design problems by temporally removing the energy constraint. Then, we propose a new mechanism design framework, called the Integrate Rounding Scheme based Maxima-in-distributional Range (IRSM), and based on that, design a new incentive mechanism for each one-shot problem. Finally, we reconsider energy constraints to design a new nonlinear online incentive mechanism by rationally combining the previously derived one-shot ones. Theoretical analyses show that our proposed nonlinear online incentive mechanism can guarantee individual rationality, truthfulness, a sound competitive ratio, and computational efficiency. We further conduct comprehensive simulations to validate the effectiveness and superiority of our proposed mechanism.

**Index Terms**—Edge computing, task offloading, online nonlinear incentive mechanism, social welfare maximization.

◆

## 1 INTRODUCE

Today, Internet of Things (IoT) technologies have been envisioned as a promising way to achieve a secure, comfortable, and convenient life. As a consequence, IoT devices are anticipated to execute a massive amount of tasks, which may cost a significant amount of computation resources, and sometimes may even exceed the local computation capacities of IoT devices [1]. To address this issue, recent prevailing edge computing technologies have been introduced, by which computation-intensive tasks can be offloaded to not only edge servers but also resourceful terminals in the proximity of IoT devices, such as ubiquitously accessible mobile users. This offloading process forms a three-tier architecture, including IoT devices, resourceful mobile users, and a central controller (such as the Base Station (BS)), which is actually an extension of the Hyrax project [2].

In practice, these resourceful mobile users may not always be willing to actively execute tasks from IoT devices if no reimbursements are provided, as they need to consume their own computation and communication resources. Besides, computation tasks from IoT devices are generated

G. Li and J. Cai are with Network Intelligence and Innovation Lab (NI2L) in the Department of Electrical and Computer Engineering, Concordia University, Montreal, QC H4B 1R6, Canada (e-mail: Gang.Li@mail.concordia.ca; Jun.Cai@concordia.ca).

X. Chen is with the VTT Technical Research Centre of Finland, Oulu 90570, Finland (e-mail: xianfu.chen@vtt.fi).

Z. Su is with the School of Cyber Science and Engineering, Xi'an Jiaotong University, Xi'an, China, and also with the School of Mechatronic Engineering and Automation, Shanghai University, Shanghai, 200444, P. R. China (email: zhousu@ieee.org).

along time, so that offloading decisions at the BS have to be made in an online manner without knowing information on possible future tasks. Moreover, it is well recognized that both communication and computation resources are limited at both the edge server and the mobile users. Therefore, in order to rationally utilize these limited resources and serve more tasks with their delay requirements, jointly optimizing communication resources, including transmission power and bandwidth for both upload link (from mobile users to the BS) and download link (from the BS to mobile users), and computation resources becomes mandatory. Most importantly, mobile users are battery-powered so that they are energy-constrained, or in other words, they have energy budgets. Without careful management, it would be possible that some of mobile users may use up their energy too fast to be available for any future participation. This may result in soaring maintenance cost as the remaining mobile users may ask for more reimbursements due to the reduction of competitions. In summary, for a practical edge computing system, an incentive mechanism with the consideration of online decision making, joint computation and communication resource allocation, and energy budget has to be designed. However, designing such an online truthful mechanism is very challenging due to the following aspects.

- Joint consideration of transmission power, bandwidth, and computation resource allocation in both upload and download links introduces nonlinear constraints, and such a joint optimization problem belongs to typical combinatorial and mixed-integer programming. In addition, as aforementioned, tasks do not show up simultaneously in reality, leading to

the design of an online algorithm. In combination, the considered problem falls into the scope of the nonlinear mixed-integer online optimization, which is extremely intractable;

- The consideration of energy budget makes offloading decision-making coupled along the time, which prevents the simple solution of treating each time slot independently. Moreover, it would be difficult and non-trivial to keep a sound performance in comparison to the corresponding offline optimal solution while well balancing mobile users' energy consumption along the time;
- Similar to offline incentive mechanism designs, mobile users' private information, such as their valuation preferences, should be submitted to the BS truthfully. More importantly, reimbursements are usually calculated based on all mobile users' private information so as to guarantee truthfulness. However, since only current information on offloaded tasks has been revealed in the online scenario, it will be more challenging to prevent mobile users from misreporting such information for designing online incentive mechanisms [3].

In the literature, existing work, such as [4]-[18], mainly focused on offline or online communication or computation resource allocations under the assumption that mobile users and edge servers were willing to provide such computation service. Even though some researches [19]-[27] have studied the ways to incentivize mobile users or edge servers, energy budgets at mobile users were usually ignored and the more complicated resource allocation problem integrating power, bandwidth and computation resources, which is common for a practical system, has not been well addressed. Moreover, most of existing work considered offline incentive mechanism designs by assuming that all tasks arrived at the system at the same time. Different from all existing work, we consider a more challenging and realistic scenario in edge computing systems. To address the aforementioned challenges, in this paper, a newly designed online truthful mechanism is proposed for task offloading. The considered system model consists of IoT devices, a central controller (such as the BS), and multiple mobile users. At the beginning of each time slot, the BS firstly collects requests of offloading tasks from IoT devices and then broadcasts them to mobile users, who will then submit their valuations and available energy to the BS. After that, the BS determines the best mobile user for each task, the transmission power and bandwidth for both upload and download transmissions, the allocated computation resource, and the corresponding payment to mobile users. This process recurs along the time, and decision making in each time slot is correlated because of the energy constraint on each mobile user. With the objective of maximizing social welfare, we first formulate an offline optimization problem and design a nonlinear online truthful mechanism based on the rule of Maximal-in-Distributional Range (MIDR) [29]. Note that even though the works [30]-[32] have designed online truthful mechanisms for crowdsensing systems, they cannot be applied in our case because their considered problems were linear.

The main contributions of this paper are summarized as follows.

- With the consideration of the energy budget at each mobile user, a nonlinear online incentive mechanism, integrating mobile user selection, communication resource allocation for both upload and download links, and computation resource allocation, is proposed for task offloading in the edge computing;
- In our proposed mechanism, a new framework for incentive (truthful) mechanism design, called Integrate Rounding Scheme based MIDR (IRSM), is proposed, which is applied to design the one-shot incentive mechanism;
- We theoretically prove that the proposed online mechanism has the properties of computation efficiency with a sound competitive ratio of $\beta(1 - \frac{1}{e})(2^\phi - 1)$, incentive compatibility and individual rationality;
- Numerical simulations have been conducted to justify our theoretical analyses and verify the effectiveness of our proposed online mechanism.

The remainder of this paper is organized as follows. In the next section, the related work is introduced. In Section 3, the system model is described, and the problem formulation is elaborated. In Section 4, the proposed online truthful mechanism framework is presented. Numerical results are presented in Section 5, followed by concluding remarks in Section 6.

## 2 RELATED WORK

Task offloading in edge computing has attracted substantial attention from both industries and academia. The existing work can be roughly classified into four categories.

**Offline resource allocation:** In edge computing systems, resource allocation is one of the most popular research topics. Resources under consideration include radio resource (such as transmission power and channels) and computation resource (CPU computation frequency). These works belong to three different categories based on their optimization objectives, even though the system architectures and constraints may be varied in different scenarios: task delay minimization [4]-[6], energy consumption minimization [7], [8], and the weighted delay and energy minimization [9], [10]. Note that these works considered the offline scenario where all offloading tasks and mobile users were all available at the decision time. In contrast, the online setting allows required information arriving along time and makes offloading decision dynamically, which is more practical in reality.

**Online resource allocation:** In order to consider a more practical situation, several researchers proposed online algorithms to allocate the resources on-the-fly. For example, authors in [11]-[12], by using Markov Decision Process (MDP) method, studied the joint optimization of radio resource scheduling and computation offloading to minimize energy consumption of all the users to process their applications, while respecting the predefined delay requirements. However, the MDP based method needs to model state spaces and decision spaces, which are related to the number of users. If the number of users is large enough, this online

Table 1
Summary and Comparison of previous and the current incentive mechanism design works in edge computing systems.

| Reference | Optimization problem | Arbitrary time length? (For objective) | Competitive Ratio? (For online solution) | Transmission power allocation | Energy budget? | Truthfulness? |
|---|---|---|---|---|---|---|
| Kiani et al. [19] | Offline nonlinear | NA | NA | ✓ | ✗ | ✓ |
| Chang et al. [22] | Offline nonlinear | NA | NA | ✗ | ✗ | ✓ |
| Mashhadi et al. [23] | Offline linear | NA | NA | ✗ | ✓ | ✓ |
| Wang et al. [24] | Online linear | ✓ | ✗ | ✗ | ✗ | ✓ |
| Li et al. [25] | Online linear | ✓ | ✓ | ✗ | ✗ | ✓ |
| Zhang et al. [27] | Online nonlinear | ✗ | ✓ | ✗ | ✓ | ✓ |
| **This work** | Online nonlinear | ✓ | ✓ | ✓ | ✓ | ✓ |

approach would be infeasible due to its high computational complexity. Different from the MDP based methods, authors in [13], [14] proposed the online radio and computational resource management algorithms based on the Lyapunov (Lyapunov) method with the objective of minimizing the long-term weighted sum of the energy consumption of the mobile devices and the edge server, subject to a task buffer stability constraint. Moreover, with the advance in machine learning techniques, in [15], [16], a model-free deep reinforcement learning-based online computation offloading approach was proposed to optimize offloading decisions and communication resource utilisation. In addition, other methods, such as regularization technique [17] and multi-armed bandit theory [18], were also used to design online algorithms for resource allocations, and network and base station selection, respectively.

**Offline incentive mechanism design:** Recently, more and more studies began to focus on network economic aspects, as edge servers and available resourceful mobile users need rewards for task execution. In offline incentive mechanism designs, authors in [19] proposed a novel three hierarchical architecture in edge computing, where an auction-based profit maximization problem was formulated to maximize the gained profit. The problem considered not only the revenue of serving virtual machine demands and the electricity cost for running the computing and network facilities, but also the revenue lost due to network delays. An offline incentive mechanism for federated learning in edge networks was designed in [20], where the edge server recruited smart devices for federated learning via a Stackelberg game, and the objective was to minimize the learning time while maximizing the accuracy level of learning model. Besides, in [21], an offline reverse auction was established to encourage vehicles to share resources for others in vehicular networks where a Vickrey-Clarke-Groves (VCG)-based mechanism design problem was formulated while satisfying the desirable economical properties of truthfulness and individual rationality. In blockchain applications, a three-layer edge computing architecture was proposed in [22] to encourage smart devices to participate in the mining process and the edge service provider to provide computational resources. Moreover, with the consideration of task delay and energy consumption constraints of mobile users, in [23], an auction was proposed to maximize the benefit of the edge server through designing two related fully-connected deep neural networks.

**Online incentive mechanism design:** In order to enable

incentive mechanisms to be applicable in dynamic settings, [24]-[27] proposed online incentive mechanism designs for edge computing systems. Authors in [24] designed an online profit maximization multi-round auction for the computation resource trading between edge clouds (sellers) and mobile devices (buyers) in a competitive environment. For making full use of idle computation resources, online incentive mechanisms for collaborative task offloading in edge computing were proposed in [25], [26], where edge server motivated idle resourceful mobile users to provide their computation resources to requesters who would like to offload tasks and participate in the system dynamically. Moreover, by considering the energy harvesting process at mobile users, an optimization problem with a long-term reward objective was formulated in [27] to investigate sustainable computation offloading in an edge computing system, and then an online incentive mechanism based on the Lyapunov method and the VCG payment was designed to solve this problem. A summary of typical online mechanism design methods was provided in [28]

However, the previous works for online truthful mechanism designs in edge computing only considered either computation resource allocation or computation resource and bandwidth allocations, while ignoring transmission power allocation at both edge server and mobile users' sides, and didn't consider the energy constraint at each mobile user, which results in the considered systems to be nonlinear and intractable. Different from all existing works, we propose a new online incentive mechanism integrating power, bandwidth, and computation resource allocation for edge computing systems. For intuitive understanding and clarification, Table 1 summarizes the differences between our proposed work and existing online incentive mechanisms in edge computing systems.

## 3 SYSTEM MODEL AND PROBLEM FORMULATION

In this section, we first describe the system under consideration, including both computation and communication models, and then formulate the interaction between the BS and mobile users as an online incentive mechanism design problem. After that, the corresponding offline optimization problem is formulated. For the notational convenience, Table 2 summarizes the major notations used in this paper.

### 3.1 Network Artechiture

Similar to [27], consider an edge computing system, as shown in Fig. 1. In the system, there are a BS, several IoT

## Table 2
### Commonly Used Notations in This Paper

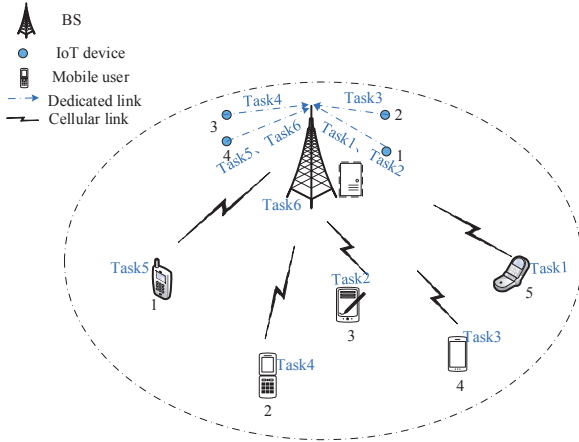| Notation | Interpretation |
|---|---|
| $\mathcal{T}$ | set of time slots |
| $\Delta T$ | time length of each time slot |
| $\mathcal{N}^t$ | Set of task at time slot $t$ |
| $\mathcal{M}^t$ | Set of available mobile users at time slot $t$ |
| $L_j^t$ | bid of mobile user $i$ at time slot $t$ |
| $c_j^t$ | submitted unit power cost by mobile user $j$ |
| $v_j^t$ | true or actual cost of $c_j^t$ |
| $c_{B,j}^t$ | unit power cost at the BS to mobile user $j$ |
| $E_j$ | total energy of mobile user $j$ |
| $L_i^t$ | set of submitted bids of the mobile user $i$ |
| $J_i^t$ | task $i$ at time slot $t$ |
| $S_i^t$ | input size of task $i$ at time slot $t$ |
| $O_i^t$ | size of returned result at time slot $t$ |
| $D_i^t$ | delay tolerance of task $i$ at time slot $t$ |
| $Q_i^t$ | total required CPU cycles of task $i$ at time slot $t$ |
| $T_{i,j}^{exe,t}$ | total execution time of task $i$ at time slot $t$ |
| $f_{i,j}^t$ | allocated CPU frequency at mobile user $j$ |
| $F_j$ | total CPU frequency at mobile user $j$ |
| $C_{i,j}^{exe,t}$ | total energy consumption of task $i$ at time slot $t$ |
| $B^U, B^D$ | total bandwidth for the downlink and uplink |
| $w_{i,j}^{D,t}, w_{i,j}^{U,t}$ | allocated downlink and uplink bandwidths for task $i$ |
| $x_{i,j}^t$ | binary decision variable |
| $e_{i,j}^t$ | total energy consumption at mobile user $j$ for task $i$ |
| $R_{i,j}^{U,t}, R_{i,j}^{D,t}$ | uplink and downlink transmission rates |
| $P_{i,j}^{U,t}, P_{i,j}^{D,t}$ | transmission power at mobile user $j$ and the BS |
| $T_{i,j}^{U,t}, T_{i,j}^{D,t}$ | uplink and downlink transmission times for task $i$ |
| $\pi_{i,j}^t$ | reimbursement for mobile user $j$ foe executing task $i$ |
| $\boldsymbol{L}^t$ | collection of all submitted bids at time slot $t$ |
| $\boldsymbol{L}_j^t$ | collection of all submitted bids by mobile user $j$ |
| $\boldsymbol{L}_{-j}^t$ | collection of bids, excluding mobile user $j$ |



Figure 1. The system model with 4 IoT devices and 5 mobile users.

devices, and some mobile users who provide computation services in the coverage of the BS. IoT devices generate tasks along the time and submit their tasks at the beginning of each time slot to the BS through dedicated links, such as WiFi or Low-power Wide Area Networking (LPWAN). By considering different computation demands for different tasks as in [33], in this paper, we consider two different types of tasks in the simulation, i.e., image decompression and decision-marking. After receiving all these tasks, the BS determines the task assignment among mobile users and itself, and corresponding computation and communication

resource allocations. After that, the BS offloads assigned tasks to the mobile users. Once the completion of tasks, results will be returned to the BS by mobile users. Note that for those unoffloaded tasks, the BS will consume its own computation resource and energy for execution[1]. Define a time-slotted structure with time slots indexed by $\mathcal{T} = \{1, 2, \cdots, T\}$, where $T$ is the total number of time slots, and there is a set of tasks $\mathcal{N}^{(t)}$ at each time slot $t$ with cardinality of $|\mathcal{N}^{(t)}| = N^{(t)}$. At the beginning of each time slot, the BS broadcasts tasks to a set of mobile users $\mathcal{M}^{(t)}$ with cardinality of $|\mathcal{M}^{(t)}| = M^{(t)}$ for execution, and mobile users connect with the BS via the cellular network. Similar to existing work in edge computing [25], [27], a quasi-static scenario is studied, where all mobile users and wireless communication configurations keep stationary in each time slot, but may change slot by slot. In addition, following [27], [34], if the task has been successfully offloaded to a mobile user in time slot $t$, its execution will be completed within the time slot period, i.e., $\Delta T$. Technically, if the size of generated task is too large, the IoT device can split the oversize task into multiple small tasks before offloading [35], [36].

Obviously, the aforementioned interactions could be modelled as an auction, where mobile users are sellers who sell their unused computational resource for monetary benefits, while the BS is the buyer. At the beginning of each time slot $t$, the BS broadcasts descriptions of $N^{(t)}$ tasks to mobile users, and then the bid from each mobile user $j$, denoted as $L_j^{(t)} = \{c_j^{(t)}, E_j\}$, is submitted to the BS. Here $c_j^{(t)}$ represents the cost per unit power consumption at mobile user $j$, and $E_j$ is the total available energy of mobile user $j$, which is submitted at the first time slot only. We denote the true or actual valuation of $c_j^{(t)}$ by $v_j^{(t)}$, which is private and only known to the mobile user $j$ [37]. After collecting all bids from mobile users, the BS determines which task should be executed by which mobile user and how much reimbursement should be given to winning mobile users.

Following [38], [39], data transmission in the download link (from the BS to mobile users) and upload link (from mobile users to the BS) are mainly considered in this paper, while the overhead for the control signalling is overlooked. This is because compared to the data size and returned results, control signalling is much smaller, and can be transmitted through dedicated channels.

### 3.2 Computation Model

We mathematically characterize each task $i$ at time slot $t$ by a tuple $J_i^{(t)} = \{S_i^{(t)}, O_i^{(t)}, D_i^{(t)}\}$, where $S_i^{(t)}$ denotes the size of the input, $O_i^{(t)}$ denotes the size of the return result, and $D_i^{(t)}$ denotes the delay tolerance. The required amount of CPU cycles $Q_i^{(t)}$ for task $i$ can be estimated as [6]

$$Q_i^{(t)} = \epsilon_i S_i^{(t)}, \tag{1}$$

where $\epsilon_i$ is the CPU cycle coefficient, which may vary from different tasks or applications, and can be obtained as shown in [40], [41]. Note that the allocated computational frequency at each mobile user $j$ is an optimization variable,

---

1. Since the offloading process between the BS and mobile users is the main focus in this paper, the computation resource allocation at the BS side is ignored.

which is represented by $f_{i,j}^{(t)}$. Then, the execution time for the allocated task at the mobile user $j$ can be calculated as

$$T_{i,j}^{exe,(t)} = \frac{Q_i^{(t)}}{f_{i,j}^{(t)}}, \ \forall j \in \mathcal{M}^{(t)}, \forall i \in \mathcal{N}^{(t)}. \quad (2)$$

Furthermore, since the computational capacity of each mobile user $j$ is ordinarily limited, the allocated computation frequency for the task should be no more than this limitation, i.e.,

$$f_{i,j}^{(t)} \le F_j, \quad \forall j \in \mathcal{M}^{(t)}, \quad (3)$$

where $F_j$ denotes the computational capacity of mobile user $j$. Based on [6] and [42], the energy consumption $E_{i,j}^{exe,(t)}$ (unit: Joule) and energy consumption cost $C_{i,j}^{exe,(t)}$ (unit: dollar) at mobile user $j$ for executing task $i$ can be respectively computed as

$$E_{i,j}^{exe,(t)} = \beta_j Q_i^{(t)} (f_{i,j}^{(t)})^2, \ \forall j \in \mathcal{M}^{(t)}, \forall i \in \mathcal{N}^{(t)}, \quad (4)$$

$$C_{i,j}^{exe,(t)} = \theta_j Q_i^{(t)}, \ \forall j \in \mathcal{M}^{(t)}, \forall i \in \mathcal{N}^{(t)}, \quad (5)$$

where $\beta_j$ is the energy consumption coefficient at mobile user $j$, and $\theta_j$ is the energy cost per CPU cycle. According to [6], [42], $\beta_j$ and $\theta_j$ can be estimated in practice.

### 3.3 Communication Model

The communication model includes download and upload links, whereby tasks can be offloaded to selected mobile users, and computation results can be fed back to the BS, respectively. Let $B^U$ and $B^D$ be the total bandwidths for the upload and download links, respectively. Moreover, a multi-channel wireless communication setting is considered, where the system's bandwidth is divided into wireless communication channels. The data transmission between each mobile user and the base station will be allocated an exclusive channel for task offloading [14], [26]. Define $x_{i,j}^{(t)}$ to be an indicator variable, which means task $i$ is allocated to mobile user $j$ in time slot $t$ if $x_{i,j}^{(t)} = 1$, and $x_{i,j}^{(t)} = 0$ otherwise.

Each mobile user can only execute at most one task at each time slot $t$, and each task can only be offloaded to at most one mobile user. Those tasks, which are not executed at the current time slot $t$, will be executed by the BS before their deadlines. For $x_{i,j}^{(t)}$, we have the following constraints

$$\sum_{i \in \mathcal{N}^{(t)}} x_{i,j}^{(t)} \le 1, \ \forall j \in \mathcal{M}^{(t)}, \forall t \in \mathcal{T}, \quad (6)$$

$$\sum_{j \in \mathcal{M}^{(t)}} x_{i,j}^{(t)} \le 1, \ \forall i \in \mathcal{N}^{(t)}, \forall t \in \mathcal{T}. \quad (7)$$

In addition, since each mobile user is energy constrained, it has the energy restriction over the time as

$$\sum_{t \in \mathcal{T}} \sum_{i \in \mathcal{N}^{(t)}} e_{i,j}^{(t)} x_{i,j}^{(t)} \le E_j, \ \forall j \in \mathcal{M}^{(t)}, \quad (8)$$

where $e_{i,j}^{(t)}$ is the total energy consumption of mobile user $i$ at time slot $t$. Note that as indicated in the constraint (8), the decision-makings are coupled along the time slots. Moreover, $e_{i,j}^{(t)}$ is a combination of upload link transmission energy and task execution energy, as shown in (12), both of which are related to optimization variables. Hence, by

introducing the energy budget on each mobile user, the considered mechanism design problem is far more challenging than traditional ones in the literature.

For the upload link, since mobile users need to transmit results back to the BS, the transmission rate with the allocated bandwidth $w_{i,j}^{U,(t)}$ can be calculated as

$$R_{i,j}^{U,(t)} = w_{i,j}^{U,(t)} \log_2(1 + \gamma_{i,j}^{U,(t)}), \quad (9)$$

where $\gamma_{i,j}^{U,(t)} = \frac{P_{i,j}^{U,(t)} G_{i,j}^{U,(t)}}{\sigma^2}$ is the signal-to-noise ratio (SNR). $P_{i,j}^{U,(t)}$ and $G_{i,j}^{U,(t)} = |h_{i,j}^{U,(t)}|^2$ are the transmission power and the channel gain between mobile user $j$ and the BS, respectively, and $\sigma^2$ represents the background noise power. $h_{i,j}^{U,(t)}$ is the channel coefficient, which is modelled as Rayleigh fading and all channel coefficients are zero-mean, circularly symmetric complex Gaussian (CSCG) random variables with variances $d^{-\frac{v}{2}}$, where $d$ is the distance between the transmitter and the receiver. Note that each mobile user has to satisfy its power constraint as

$$\sum_{i \in \mathcal{N}^{(t)}} P_{i,j}^{U,(t)} x_{i,j}^{(t)} \le P_j^{max}, \ \forall j \in \mathcal{M}^{(t)}, \forall t \in \mathcal{T}, \quad (10)$$

where $P_j^{max}$ is the maximal transmission power of mobile user $j$. Then, the upload link transmission time for task $i$ from mobile user $j$ can be calculated as

$$T_{i,j}^{U,(t)} = \frac{O_i^{(t)}}{R_{i,j}^{U,(t)}}. \quad (11)$$

Combine constraints (4) and (11), we have

$$e_{i,j}^{(t)} = P_{i,j}^{U,(t)} \times T_{i,j}^{U,(t)} + E_{i,j}^{exe,(t)}. \quad (12)$$

As for the download link transmission, the BS needs to transmit allocated tasks to selected mobile users. The transmission rate between the mobile user $j$ and the BS for task $i$ can be calculated as

$$R_{i,j}^{D,(t)} = w_{i,j}^{D,(t)} \log_2(1 + \gamma_{i,j}^{D,(t)}), \quad (13)$$

where $w_{i,j}^{D,(t)}$ is the allocated bandwidth for the download link channel, and $\gamma_{i,j}^{D,(t)} = \frac{P_{i,j}^{D,(t)} G_{i,j}^{D,(t)}}{\sigma^2}$ represents the SNR. $P_{i,j}^{D,(t)}$ and $G_{i,j}^{D,(t)} = |h_{i,j}^{D,(t)}|^2$ are the allocated transmission power and download link channel gain for task $i$ to mobile user $j$, respectively. $h_{i,j}^{D,(t)}$ is the download link channel coefficient, which is also modelled as Rayleigh fading. Similarly, the download link transmission time can be expressed as

$$T_{i,j}^{D,(t)} = \frac{S_i^{(t)}}{R_{i,j}^{D,(t)}}. \quad (14)$$

Since in wireless communication networks, the BS commonly has a transmission power limit $P_B^{max}$, we have

$$\sum_{i \in \mathcal{N}^{(t)}} \sum_{j \in \mathcal{M}^{(t)}} P_{i,j}^{D,(t)} x_{i,j}^{(t)} \le P_B^{max}, \ \forall t \in \mathcal{T}. \quad (15)$$

With the consideration of the limited bandwidths for both the download and upload links, the following conditions should be imposed on the bandwidth allocation as

$$\sum_{i \in \mathcal{N}^{(t)}} \sum_{j \in \mathcal{M}^{(t)}} w_{i,j}^{D,(t)} x_{i,j}^{(t)} \le B^D, \ \forall t \in \mathcal{T}, \quad (16)$$

$$\sum_{i \in \mathcal{N}^{(t)}} \sum_{j \in \mathcal{M}^{(t)}} w_{i,j}^{U,(t)} x_{i,j}^{(t)} \le B^U, \ \forall t \in \mathcal{T}. \quad (17)$$

In order to meet the deadline requirement of each of-floaded task, each computation task $J_i^{(t)}$ with input data size $S_i^{(t)}$ and output data size $O_i^{(t)}$ has to be transmitted and executed within $D_i^{(t)}$. This means we have the following constrain.

$$\sum_{j \in \mathcal{M}^{(t)}} (T_{i,j}^{D,(t)} + T_{i,j}^{exe,(t)} + T_{i,j}^{U,(t)})x_{i,j}^{(t)} \leq \min\{D_i^{(t)}, \Delta T\},$$

$$\forall i \in \mathcal{N}^{(t)}, \ \forall t \in \mathcal{T}. \quad (18)$$

Note that, as shown in (18), the total delay for each task is the summation of download transmission time, upload transmission time, and execution time. Moreover, since the download or upload transmission time is actually a ratio between the data size and the transmission rate, i.e., $R_{i,j}^{D,(t)}$ or $R_{i,j}^{U,(t)}$, the constraint (18) becomes nonlinear with respect to both transmission power and bandwidth. This requires us to consider a nonlinear online incentive mechanism design. However, designing such an online incentive mechanism is extremely challenging due to the nonlinearity and the multi-dimensional allocation outcome.

### 3.4 Utility of the BS

The utility of the BS consists of the benefit $r_{i,j}^{(t)}$ through offloading tasks to mobile users, the execution cost for the BS itself, the cost for the BS to transmit tasks to mobile users, and the rewards to mobile users, i.e., $\pi_{i,j}^{(t)}$. We mathematically formulate the utility of the BS as follows

$$U_B(\boldsymbol{x}^{(t)}) = \sum_{i \in \mathcal{N}^{(t)}} \sum_{j \in \mathcal{M}^{(t)}} r_{i,j}^{(t)} x_{i,j}^{(t)} - \sum_{i \in \mathcal{N}^{(t)}} (1 - \sum_{j \in \mathcal{M}^{(t)}} x_{i,j}^{(t)})\phi_i^{(t)}$$

$$- \sum_{i \in \mathcal{N}^{(t)}} \sum_{j \in \mathcal{M}^{(t)}} P_{i,j}^{D,(t)} c_B^{(t)} x_{i,j}^{(t)} - \sum_{i \in \mathcal{N}^{(t)}} \sum_{j \in \mathcal{M}^{(t)}} \pi_{i,j}^{(t)}, \quad (19)$$

where $\boldsymbol{x}^{(t)} = \{x_{1,1}^{(t)}, x_{1,2}^{(t)}, \cdots, x_{N^{(t)},M^{(t)}}^{(t)}\}$, $\phi_i^{(t)} = \theta_B Q_i^{(t)}$ is the energy consumption cost for the BS, and $c_B^{(t)}$ and $\theta_B$ are the transmission cost per unit power and energy cost per CPU cycle, respectively. $\pi_{i,j}^{(t)}$ denotes the reimbursement from the BS to mobile user $j$ for task $i$. The BS can gain a benefit $r_{i,j}^{(t)}$ because the offloading process can, to large extent, preserve its own computation resource so that the BS will own enough resource for other intensive and complex applications.

### 3.5 Utility of Mobile User

The utility of any mobile user $j$ at any time slot $t$ consists of the payment from BS and the total cost related to executing task $i$. Thus, the utility of mobile user $j$ can be formulated as

$$U_j(\boldsymbol{L}^{(t)}) = \sum_{i \in \mathcal{N}^{(t)}} (\pi_{i,j}^{(t)} - (P_{i,j}^{U,(t)} c_j^{(t)} + \theta_j Q_i^{(t)})x_{i,j}^{(t)}),$$

$$\forall j \in \mathcal{M}^{(t)}, \ \forall t \in \mathcal{T}, \quad (20)$$

where $\boldsymbol{L}^{(t)} = \{L_1^{(t)}, L_2^{(t)}, \cdots, L_{M^{(t)}}^{(t)}\}$ is a collection of all submitted bids at time slot $t$. Intuitively, since mobile users are intelligent and selfish, and their submitted information is private and unknown to the BS, they may submit their biding information strategically in order to earn more benefits in the competition. To this end, an incentive mechanism should be designed to force mobile users to bid their private information truthfully.

### 3.6 Problem Formulation

Before formulating our problem mathematically, we first introduce several properties that should be satisfied as follows.

- Incentive Compatibility (IC) guarantees that no mobile users can gain more benefits by misreporting their private information. This property can be mathematically expressed as

$$U_j(L_j^{(t)}, \boldsymbol{L}_{-j}^{(t)}) \geq U_j(\acute{L}_j^{(t)}, \boldsymbol{L}_{-j}^{(t)}), \ \forall j \in \mathcal{M}^{(t)}, \ \forall t \in \mathcal{T}, \quad (21)$$

  where $\boldsymbol{L}_{-j}^{(t)}$ is the collection of bids excluding $L_j^{(t)}$, and $\acute{L}_j^{(t)}$ is a vector of potential false bids from mobile user $j$. Sometimes, IC can also be called as truthfulness.

- Individual Rationality (IR) ensures utilities of all mobile users are no less than zero, which can be formulated as

$$U_j(\boldsymbol{L}_j^{(t)}) \geq 0, \ \forall j \in \mathcal{M}^{(t)}, \ \forall t \in \mathcal{T}. \quad (22)$$

  In reality, this property is of great importance as it attracts mobile users to participate the campaign.

- Competitive Ratio (CR) is defined as the ratio of the calculated online solutions over the optimal offline ones. This metric is vital because it evaluates the performance of the designed online incentive mechanism. Note that in this paper, the upper bound of CR is one, and the larger CR is, the better performance becomes.

- Computational Efficiency (CE) is a metric that measures whether the designed online incentive mechanism can run in a polynomial time.

The BS has to jointly determine power allocations at both the BS and mobile users $(P_{i,j}^{D,(t)}, P_{i,j}^{U,(t)})$, upload and download link bandwidths $(w_{i,j}^{D,(t)}, w_{i,j}^{U,(t)})$, computation frequency $(f_{i,j}^{(t)})$, reimbursement to the mobile user $(\pi_{i,j}^{(t)})$, and task assignment $(x_{i,j}^{(t)})$ on-the-fly. The objective is to maximize the social welfare $SW_{ob}$ of the system, which is defined as the summation of utilities from the BS and all mobile users, and can be calculated as

$$SW_{ob} = \sum_{t \in \mathcal{T}} U_B(\boldsymbol{x}^{(t)}) + \sum_{j \in \mathcal{M}^{(t)}} U_j(\boldsymbol{L}^{(t)}). \quad (23)$$

After some algebraic manipulations, the objective $SW_{ob}$ can be equivalently rewritten as

$$SW_{ob} = \sum_{t \in \mathcal{T}} \sum_{i \in \mathcal{N}^{(t)}} \sum_{j \in \mathcal{M}^{(t)}} a_{i,j}^{(t)} x_{i,j}^{(t)} - \sum_{t \in \mathcal{T}} \sum_{i \in \mathcal{N}^{(t)}} \sum_{j \in \mathcal{M}^{(t)}} P_{i,j}^{D,(t)} c_B^{(t)} x_{i,j}^{(t)}$$

$$- \sum_{t \in \mathcal{T}} \sum_{i \in \mathcal{N}^{(t)}} \sum_{j \in \mathcal{M}^{(t)}} P_{i,j}^{U,(t)} c_j^{(t)} x_{i,j}^{(t)} - \sum_{t \in \mathcal{T}} \sum_{i \in \mathcal{N}^{(t)}} \phi_i^{(t)}, \quad (24)$$

where $a_{i,j}^{(t)} = r_{i,j}^{(t)} + \phi_i^{(t)} - \theta_j Q_i^{(t)}$. Note that since the last term in (24) is a constant, which cannot affect the optimality of the optimization problem, we will ignore this term in the following calculations and analysis.

Therefore, if all required information for all time slots can be obtained, the offline joint computation and communication resource allocation problem can be formulated as

$[\mathcal{P}1]$ :

$$[\boldsymbol{X}, \boldsymbol{\Pi}] = \arg\max_{\boldsymbol{X}, \boldsymbol{\Pi}} SW_{ob}$$

$$s.t. \quad (3), (6) - (8), (10), (15) - (18), (21) - (22),$$
$$x_{i,j}^{(t)} \in \{0,1\}, w_{i,j}^{D,(t)} \geq 0, w_{i,j}^{U,(t)} \geq 0, P_{i,j}^{D,(t)} \geq 0, P_{i,j}^{U,(t)} \geq 0,$$
$$f_{i,j}^{(t)} \geq 0, \pi_{i,j}^{(t)} \geq 0, \, \forall i \in \mathcal{N}^{(t)}, \forall j \in \mathcal{M}^{(t)}, \forall t \in \mathcal{T},$$

where $\boldsymbol{X} = \{x_{i,j}^{(t)}, w_{i,j}^{D,(t)}, w_{i,j}^{U,(t)}, P_{i,j}^{D,(t)}, P_{i,j}^{U,(t)}, f_{i,j}^{(t)}\}$, and $\boldsymbol{\Pi} = \{\pi_{i,j}^{(t)}\}$, $\{i \in \mathcal{N}^{(t)}, j \in \mathcal{M}^{(t)}, t \in \mathcal{T}\}$ are decision variables. Obviously, it is very difficult, if not impossible, to solve $[\mathcal{P}1]$ in an online manner with a good CR because i) this problem is a typical mixed integer nonconvex optimization problem, which is NP-hard (refer to **Lemma** 1); ii) the constraint (8) in $[\mathcal{P}1]$ couples all time slots so that we cannot simply solve the problem for each time slot independently and unreasonable decisions in current time slot may affect future results; and iii) due to the consideration of transmission power for both download and upload links, constraints (8) and (18) become nonlinear and tightly coupled with each other, which further complicates the problem. Note that the above online mechanism design problem may be potentially solved by the MDP or the Lyapunov method. However, these two methods actually cannot be applied here because i) a *prior* distribution on biding information of mobile users, necessary for the MDP method, is not known in our case; ii) in the Lyapunov method, the objective should be an infinite time average form, while in this paper, we consider any arbitrary time span. Moreover, only by carefully meeting certain requirements can the convergence of Lyapunov method be obtained. To this end, in the following sections, a novel online incentive mechanism with a sound CR is proposed.

**Lemma 1.** *The offline incentive mechanism design problem $[\mathcal{P}1]$ is NP-hard.*

*Proof.* Note that the problem $[\mathcal{P}1]$ is nonlinear because both the task delay requirement and the energy budget constraints, i.e., equations (8) and (18), are nonlinear. However, if constraints, (8), (18), (21), and (22) are removed from the problem $[\mathcal{P}1]$, the remaining becomes a mixed integer linear optimization programming problem, which has been proved to be NP-hard [43], [44]. Since a reduced problem of $[\mathcal{P}1]$ is NP-hard, the original problem is NP-hard as well. ∎

## 4 ONLINE INCENTIVE MECHANISM DESIGN

In this section, we will present our online solution for the problem $[\mathcal{P}1]$, which consists of two steps. In the first step, by temporarily ignoring the energy constraint on each mobile user, an incentive mechanism, called o͟ne-shot t͟ruthful m͟echanism (OTM), will be designed for each single time slot. After that, we will rationally combine these independent single time slot solutions to design an online truthful mechanism.

### 4.1 One-shot Truthful Mechanism Design

In this subsection, we consider a single time slot and temporarily remove the energy constraint for each mobile user,

i.e., the constraint (8). Hence, the original problem $[\mathcal{P}1]$ can be transformed into $[\mathcal{P}2]$ as follows. For the notional simplification, the superscript $t$ in variables and sets are omitted in this subsection.

$[\mathcal{P}2]$ :

$$\max_{\boldsymbol{X}, \boldsymbol{\Pi}} \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} a_{i,j} x_{i,j} - \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} P_{i,j}^D c_B x_{i,j} - \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} P_{i,j}^U c_j x_{i,j}$$

$$s.t. \sum_{i \in \mathcal{N}} x_{i,j} \leq 1, \, \forall j \in \mathcal{M}, \tag{25}$$

$$\sum_{j \in \mathcal{M}} x_{i,j} \leq 1, \, \forall i \in \mathcal{N}, \tag{26}$$

$$\sum_{i \in \mathcal{N}} P_{i,j}^U x_{i,j} \leq P_j^{max}, \, \forall j \in \mathcal{M}, \tag{27}$$

$$\sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} P_{i,j}^D x_{i,j} \leq P_B^{max}, \tag{28}$$

$$\sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} w_{i,j}^D x_{i,j} \leq B^D, \tag{29}$$

$$\sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} w_{i,j}^U x_{i,j} \leq B^U, \tag{30}$$

$$\sum_{j \in \mathcal{M}} (T_{i,j}^D + T_{i,j}^{exe} + T_{i,j}^U) x_{i,j} \leq \min\{D_i, \Delta T\}, \forall i \in \mathcal{N} \tag{31}$$

$$T_{i,j}^D = \frac{S_i}{w_{i,j}^D x_{i,j} \log_2 (1 + \frac{P_{i,j}^D x_{i,j} G_{i,j}^D}{\sigma^2})},$$

$$T_{i,j}^U = \frac{O_i}{w_{i,j}^U x_{i,j} \log_2 (1 + \frac{P_{i,j}^U x_{i,j} G_{i,j}^U}{\sigma^2})}, \, T_{i,j}^{exe} = \frac{Q_i}{f_{i,j}},$$

$$f_{i,j} \leq F_j, \, \forall j \in \mathcal{M}, \tag{32}$$

$$U_j(\boldsymbol{L}) \geq 0, \, \forall j \in \mathcal{M}, \tag{33}$$

$$U_j(L_j, \boldsymbol{L}_{-j}) \geq U_j(\acute{L}_j, \boldsymbol{L}_{-j}), \, \forall j \in \mathcal{M}. \tag{34}$$

Moreover, since the reward or the payment, i.e., $\pi_{i,j}^{(t)}$, is not in the objective function of $[\mathcal{P}2]$, constraints (33) and (34), which respectively represent IR and IC, can be further removed. Note that this manipulation doesn't affect the optimality of $[\mathcal{P}2]$, and we will later reconsider constraints (33) and (34) by designing a payment rule. For clarity, we rewrite the newly formed optimization problem, i.e., $[\mathcal{RP}]$ as follows, which is termed as resource allocation problem in this paper.

$[\mathcal{RP}]$ :

$$\max_{\boldsymbol{X}} \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} a_{i,j} x_{i,j} - \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} P_{i,j}^D c_B x_{i,j} - \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} P_{i,j}^U c_j x_{i,j}$$

$$s.t. \quad (25) - (32)$$

Note that the problem $[\mathcal{RP}]$ is still NP-hard due to the nonconvexity and mixed integer optimization.

Generally, randomized rounding algorithms [45] are a potential solutions for problem $[\mathcal{RP}]$, which commonly consists of a relaxation algorithm $\mathcal{A}$ and a rounding algorithm $\mathcal{D}(\dot{\boldsymbol{X}})$. $\mathcal{A}$ is used to obtain the fractional solution of $[\mathcal{RP}]$, i.e., $\dot{\boldsymbol{X}}$, while $\mathcal{D}(\dot{\boldsymbol{X}})$ maps $\dot{\boldsymbol{X}}$ to integer solutions based on a predefined rounding scheme, such as poisson rounding scheme [46]. Unfortunately, as shown in [47], the outcomes from randomized rounding algorithms could hardly meet constraints of IR and IC. Instead, in this paper, by considering the extreme complexity of $[\mathcal{RP}]$ and motivated by

the Maximal-in-Distributional Range (MIDR) rule [29], we propose a new method, called integrate rounding scheme based MIDR (IRSM), for the one-shot truthful mechanism.

The basic idea of our proposed IRSM can be explained as follows. In IRSM, we integrate the integer solution $\boldsymbol{X}$ and the rounding scheme $\mathcal{D}(\dot{\boldsymbol{X}})$ into the objective function $f(\boldsymbol{X}, \boldsymbol{L})$, and then find a fractional solution $\dot{\boldsymbol{X}}$ that optimizes the expected objective function $\mathbb{E}_{\boldsymbol{X} \sim \mathcal{D}(\dot{\boldsymbol{X}})}\{f(\boldsymbol{X}, \boldsymbol{L})\}$ among all feasible fractional solutions. After that, the rounding scheme, i.e., $\mathcal{D}(\dot{\boldsymbol{X}})$, is used to obtain the integer solutions $\boldsymbol{X}$ from fractional solutions $\dot{\boldsymbol{X}}$. Although this optimization problem is usually intractable due to the embedding of the rounding function into the objective function, we can still manage it in this paper. Note that there exists a natural distinction between our proposed IRSM and the traditional randomized rounding algorithm. In fact, our IRSM is to integrate the rounding scheme in the objective and directly optimize the integer variable $\boldsymbol{X}$, rather than the relaxed original problem, which is optimized in the traditional randomized rounding algorithm. Algorithm 1 presents the steps of our proposed IRSM algorithm.

---

**Algorithm 1:** The framework of our proposed IRSM algorithm.

---

**Input:** Submitted biddings $\boldsymbol{L}$.
**Output:** Feasible solution $\boldsymbol{X}$.
1 Obtain the optimal fractional solution $\dot{\boldsymbol{X}}^*$ by maximizing $\mathbb{E}_{\boldsymbol{X} \sim \mathcal{D}(\dot{\boldsymbol{X}})}\{f(\boldsymbol{X}, \boldsymbol{L})\}$;
2 Rounding the final solution $\boldsymbol{X}$ based on the rounding scheme $\mathcal{D}(\dot{\boldsymbol{X}}^*)$.

---

Note that for most rounding schemes in the literature, the expected maximization problem in the Algorithm 1 cannot be solved in polynomial time [46]. Besides, in order to apply Algorithm 1 to design OTM, this expected maximization problem should be solved optimally. To this end, in this paper, the widely used rounding scheme, i.e., poisson rounding scheme, is applied in our analyses, which can well balance the computational complexity and the achievable performance. We will prove later that by applying this rounding scheme, the designed OTM can achieve an approximation ratio of $1 - \frac{1}{e}$ with polynomial running time, where $e$, known as the Euler's number, is a mathematical constant. Note that the proposed framework, i.e., **Algorithm** 1, can be applied to other rounding schemes.

Therefore, we reformulate $[\mathcal{RP}]$ by taking the following two operations:

- Relax the integer variable, i.e., $x_{i,j}$, to be a real variable, i.e., $\dot{x}_{i,j}$, between 0 and 1;
- Take the expectation on the objective function in the $[\mathcal{RP}]$ based on the poisson rounding scheme $\mathcal{D}(\dot{\boldsymbol{X}})$.

Then, the problem $[\mathcal{RP}]$ can be reformulated

$[\mathcal{RPO}]:$
$$\max_{\dot{\boldsymbol{X}}} \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} \mathbb{E}_{\boldsymbol{X} \sim \mathcal{D}(\dot{\boldsymbol{X}})}\{(a_{i,j} - P_{i,j}^D \varpi_{B,j} - P_{i,j}^U \varpi_j)x_{i,j}\}$$
$$s.t. \quad (25) - (32).$$

where $\dot{\boldsymbol{X}}$ is the set of fractional variables of $[\mathcal{RPO}]$, and $\varpi_{B,j}$ and $\varpi_j$ are virtual costs related to $c_B$ and $c_{i,j}$, respectively, which will be further explained in subsection 4.2. Define $\boldsymbol{P}^D$ and $\boldsymbol{P}^U$ as the sets of $P_{i,j}^D$ and $P_{i,j}^U$, respectively, and let $\tilde{P}_{i,j}^U = P_{i,j}^U x_{i,j}$, $\tilde{P}_{i,j}^D = P_{i,j}^D x_{i,j}$, $\tilde{w}_{i,j}^U = w_{i,j}^U x_{i,j}$, and $\tilde{w}_{i,j}^D = w_{i,j}^D x_{i,j}$. By applying the poisson rounding scheme in $[\mathcal{RPO}]$, we have

$[\mathcal{ERPO}]:$
$$\max_{\dot{\boldsymbol{X}}} \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} (a_{i,j} - \tilde{P}_{i,j}^D \varpi_{B,j} - \tilde{P}_{i,j}^U \varpi_j)(1 - e^{-x_{i,j}})$$
$$s.t. \quad (25), \quad (26), \tag{35}$$
$$\sum_{i \in \mathcal{N}} \tilde{P}_{i,j}^U \le P_j^{max}, \ \forall j \in \mathcal{M}, \ \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} \tilde{P}_{i,j}^D \le P_B^{max}, \tag{36}$$
$$\sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} \tilde{w}_{i,j}^D \le B^D, \quad \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} \tilde{w}_{i,j}^U \le B^U, \tag{37}$$
$$\sum_{j \in \mathcal{M}} (\tilde{T}_{i,j}^D + T_{i,j}^{exe} + \tilde{T}_{i,j}^U)x_{i,j} \le \min\{D_i, \Delta T\}, \tag{38}$$
$$\tilde{T}_{i,j}^D = \frac{S_i}{\tilde{w}_{i,j}^D \log_2(1 + \frac{\tilde{P}_{i,j}^{D,t} G_{i,j}^D}{\sigma^2})}, \ T_{i,j}^{exe} = \frac{Q_i}{f_{i,j}}, \ f_{i,j} \le F_j, \tag{39}$$
$$\tilde{T}_{i,j}^U = \frac{O_i}{\tilde{w}_{i,j}^U \log_2(1 + \frac{\tilde{P}_{i,j}^U G_{i,j}^U}{\sigma^2})}. \tag{40}$$

Obviously, even though $x_{i,j}$ is fractional, it is still challenging to solve $[\mathcal{ERPO}]$ because both the objective and the constraint (38) are non-convex. To address this issue, we introduce the following substitutes:

$$x_{i,j} = -\chi_{i,j}^3, \tag{41}$$
$$\ln(1 + \frac{\tilde{P}_{i,j}^{D,t} G_{i,j}^D}{\sigma^2}) = \alpha_{i,j}^D, \tag{42}$$
$$\ln(1 + \frac{\tilde{P}_{i,j}^{U,t} G_{i,j}^U}{\sigma^2}) = \alpha_{i,j}^U, \tag{43}$$
$$e^{\alpha_{i,j}^D + \chi_{i,j}^3} = z_{i,j}^D, \tag{44}$$
$$e^{\alpha_{i,j}^U + \chi_{i,j}^3} = z_{i,j}^U. \tag{45}$$

Then, the optimization problem $[\mathcal{ERPO}]$ can be reformulated as

$[\mathcal{ERPO}1]:$
$$\max \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} a_{i,j}(1 - e^{\chi_{i,j}^3}) - \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} b_{i,j}(e^{\alpha_{i,j}^D} + e^{\chi_{i,j}^3} - z_{i,j}^D - 1)$$
$$- \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} d_{i,j}(e^{\alpha_{i,j}^U} + e^{\chi_{i,j}^3} - z_{i,j}^U - 1)$$
$$s.t. -\sum_{i \in \mathcal{N}} \chi_{i,j}^3 \le 1, \ \forall j \in \mathcal{M}, \quad -\sum_{j \in \mathcal{M}} \chi_{i,j}^3 \le 1, \ \forall i \in \mathcal{N}, \tag{46}$$
$$\sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} \tilde{w}_{i,j}^D \le B^D, \quad \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} \tilde{w}_{i,j}^U \le B^U, \tag{47}$$

$$\sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} \frac{\sigma^2}{G_{i,j}^D}(e^{\alpha_{i,j}^D} - 1) \leq P_B^{max}, \tag{48}$$

$$\sum_{i \in \mathcal{N}} \frac{\sigma^2}{G_{i,j}^U}(e^{\alpha_{i,j}^U} - 1) \leq P_j^{max}, \ \forall j \in \mathcal{M}, \tag{49}$$

$$\sum_{j \in \mathcal{M}} \left( \frac{\chi_{i,j}^3 S_i}{\tilde{w}_{i,j}^D \alpha_{i,j}^D \log_2 e} + \frac{\chi_{i,j}^3 O_i}{\tilde{w}_{i,j}^U \alpha_{i,j}^U \log_2 e} + \frac{\chi_{i,j}^3 Q_i}{f_{i,j}} \right) \geq$$
$$- \min\{D_i, \Delta T\}, \tag{50}$$

$$f_{i,j} \leq F_j, \ \forall j \in \mathcal{M}, \tag{51}$$

$$(44) - (45),$$

$$\chi_{i,j} \in [-1,0], \alpha_{i,j}^D \geq 0, \alpha_{i,j}^U \geq 0, \tilde{w}_{i,j}^D \geq 0, \tilde{w}_{i,j}^U \geq 0,$$
$$z_{i,j}^D \geq 0, z_{i,j}^U \geq 0, f_{i,j} \geq 0,$$

where $b_{i,j} = \frac{\sigma^2 \varpi_{B,j}}{G_{i,j}^D}$ and $d_{i,j} = \frac{\sigma^2 \varpi_j}{G_{i,j}^U}$.

**Lemma 2.** *The transformed optimization problem $[\mathcal{ERPO}1]$ is a convex optimization problem.*

*Proof.* It can be easily derived that constraints (46)-(49) and (51) are all convex. In the following, we mainly prove that the objective function is concave, and constraints (44), (45), and (50) are all convex.

- Concavity of the objective function.
  Since the function $e^{kx}$ is convex with parameter $k$, $1 - e^{kx}$ is concave, indicating that the first term of objective function is concave. The second term is concave due to the fact that the summation of two convex functions and a linear function is still convex. Similarly, the last term is concave too. In summary, the objective function, which is the summation of three concave functions, is concave.
- Convexity of the constraint (50).
  Note that the first and the second terms in the constraint (50) have the same mathematical form of $f(x,y,z) = \frac{x^3}{yz}$, while the last term in the constraint (50) has the form of $g(x,y) = \frac{x^3}{y}$. Thus, we only need to prove $f(x,y,z)$ and $g(x,y)$ are convex. Calculate the second partial derivation of $f(x,y,z)$ with respect to $x,y,z$ as

$$\Delta^2(x,y,z) = \begin{bmatrix} \frac{6x}{yz} & -3\frac{x^2y^2}{z} & -3\frac{x^2z^2}{y} \\ -3\frac{x^2y^2}{z} & 2\frac{x^3y^3}{z} & x^3(yz)^2 \\ -3\frac{x^2z^2}{y} & x^3(yz)^2 & 2\frac{x^3z^3}{y} \end{bmatrix} \tag{52}$$

  Obviously, the first subdeterminant orders with respective to $x,y,z$ are all less than zero. The second and the third subdeterminant orders of $\Delta^2(x,y,z)$ can be respectively calculated as

$$Det_2(\Delta^2(x,y,z)) = 3(xy)^4 z^2 \geq 0,$$
$$Det_3(\Delta^2(x,y,z)) = 0.$$

  Therefore, $f(x,y,z)$ is concave. Since $g(x,y)$ is a part of $f(x,y,z)$, it is easy to prove that its first subdeterminant order is less than zero, while its second subdeterminant order is no less than zero, which means $g(x,y)$ is concave. Since the summation of concave functions is still concave, the constraint (50) is concave.
- Convexity of constraints (44) and (45).

After some simple manipulations, constraints (44) and (45) can be changed to $\alpha_{i,j}^D + \chi_{i,j}^3 - \ln z_{i,j}^D = 0$ and $\alpha_{i,j}^U + \chi_{i,j}^3 - \ln z_{i,j}^U = 0$, respectively. Thus, both (44) and (45) follow the same mathematical form of $f(x,y,z) = x^3 + y + \ln z^{-1}$. It is easy to prove that $f(x,y,z)$ is convex because a summation of two convex functions, i.e., $x^3$ and $\ln z^{-1}$, and a linear function is also convex. Therefore, both constraints (44) and (45) are convex.

Since the objective is concave and all constraints are convex, the optimization problem $[\mathcal{ERPO}1]$ is convex. This completes the proof. ∎

Since the optimization problem $[\mathcal{ERPO}1]$ is convex, we can solve it optimally by the Lagrange dual technique [50]. Denote the outcomes of problem $[\mathcal{EPRO}]$ as $\dot{X}^* = \{x_{i,j}^*, w_{i,j}^{D*}, w_{i,j}^{U*}, P_{i,j}^{D*}, P_{i,j}^{U*}, f_{i,j}^*\}$, and the final result is determined as $X = \{x_{i,j}, w_{i,j}^D, w_{i,j}^U, P_{i,j}^D, P_{i,j}^U, f_{i,j}\}$ by adopting poisson rounding scheme.

After solving the resource allocation problem, we move to design a payment scheme $\pi_{i,j}$ to satisfy constraints (33) and (34). Note that unlike existing work [30], [51] where allocation result was single-dimensional so that the Myerson lemma can be directly applied. In this paper, the allocation outcome is six-dimensional, including upload and download link power controls, upload and download link bandwidth allocations, task assignment, and computation resource allocation. Thus, we design a new payment scheme based on the framework of Vickrey-Clarke-Groves (VCG) mechanism. The designed reward or the payment scheme for mobile user $j$ for the execution of task $i$ is

$$\pi_{i,j} = \frac{\sum\limits_{i \in \mathcal{N}}(P_{i,j}^U \varpi_j + \theta_j Q_i) x_{i,j}}{\sum\limits_{i \in \mathcal{N}}(P_{i,j}^{U*} \varpi_j + \theta_j Q_i)(1 - e^{-x_{i,j}^*})} \pi_{i,j}^f, \ \forall j \in \mathcal{U}^{sel},$$

where $\mathcal{U}^{sel}$ is the set of selected mobile users, and $\pi_{i,j}^f$ is defined as the fractional payment, which can be calculated as

$$\pi_{i,j}^f = \sum_{i \in \mathcal{N}} \sum_{j' \in \mathcal{M}/\{j\}} (a_{i,j'} - P_{i,j'}^{D*}\varpi_{B,j'} - P_{i,j'}^{U*}\varpi_{j'})(1 - e^{-x_{i,j'}^*})$$
$$- \max_{X} \sum_{i \in \mathcal{N}} \sum_{j' \in \mathcal{M}/\{j\}} (a_{i,j'} - P_{i,j'}^D\varpi_{B,j'} - P_{i,j'}^U\varpi_{j'})(1 - e^{-x_{i,j'}}).$$

Till now, the whole design process of OTM for solving $[\mathcal{P}2]$ has been completed. For better understanding, we summarize the whole procedures of OTM in **Algorithm 2**.

**Lemma 3.** *The approximation ratio of the proposed **Algorithm 2** is $1 - \frac{1}{e}$ in expectation.*

*Proof.* We define $SW_{ob}$ as the objective value by the **Algorithm 2**, and $x_{i,j}$, $P_{i,j}^D$, and $P_{i,j}^U$ are corresponding solutions.

---

**Algorithm 2:** One-shot truthful mechanism by proposed IRSM.

1 **Initialization**;
2 $\mathcal{U}^{sel} = \emptyset$;
3 $x_{i,j} = 0, \; \forall\, i \in \mathcal{N}, \; \forall\, j \in \mathcal{M}$;
4 Solve the convex problem $[\mathcal{ERPO}1]$ to obtain fraction solutions $x_{i,j}^*, w_{i,j}^{D*}, w_{i,j}^{U*}, P_{i,j}^{D*}, P_{i,j}^{U*}, f_{i,j}^*$;
5 **while** $i \in \mathcal{N}$ **do**
6    Draw $d_i$ uniformly at random from $[0, 1]$;
7    **if** $\sum\limits_{j \in (\mathcal{M}/\mathcal{U})} (1 - e^{-x_{i,j}^*}) \geq d_i$ **then**
8      Let $j^*$ be the minimum index to satisfy that $\sum\limits_{j \leq j^*} (1 - e^{-x_{i,j}^*}) \geq d_i$;
9      $\mathcal{U}^{sel} = \mathcal{U}^{sel} \cup j^*$;
10      $x_{i,j^*} = 1$ ;
11 Make the payment to selected mobile user $j$ as

$$\pi_{i,j} = \frac{\sum\limits_{i \in \mathcal{N}} (P_{i,j}^U \varpi_j + \theta_j Q_i) x_{i,j}}{\sum\limits_{i \in \mathcal{N}} (P_{i,j}^{U*} \varpi_j + \theta_j Q_i)(1 - e^{-x_{i,j}^*})} \pi_{i,j}^f;$$

12 **Output**;
13 $x_{i,j}, w_{i,j}^D, w_{i,j}^U, P_{i,j}^D, P_{i,j}^U, f_{i,j}, \pi_{i,j}$ ;

---

Moreover, let $SW_{ob}^*$ be the optimal objective value of the problem $[\mathcal{RP}]$. Then, we have

$$\mathbb{E}\{SW_{ob}\} = \mathbb{E}\Big\{ \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} (a_{i,j} - P_{i,j}^D \varpi_{B,j} - P_{i,j}^U \varpi_j) x_{i,j} \Big\}$$

$$= \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} (a_{i,j} - P_{i,j}^D \varpi_{B,j} - P_{i,j}^U \varpi_j)(1 - e^{-x_{i,j}})$$

$$\geq \Big( \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} (a_{i,j} - P_{i,j}^{D*} \varpi_{B,j} - P_{i,j}^{U*} \varpi_j)(1 - e^{-x_{i,j}^*}) \Big)$$

$$\geq (1 - \frac{1}{e}) \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} (a_{i,j} - P_{i,j}^{D*} \varpi_{B,j} - P_{i,j}^{U*} \varpi_j) x_{i,j}^*$$

$$\geq (1 - \frac{1}{e}) SW_{ob}^*.$$

This completes the proof. ∎

**Lemma 4.** *The proposed mechanism (OTM) is individually rationale in expectation.*

*Proof.* We have

$$\mathbb{E}\{U_j(\boldsymbol{L})\} = (1 - e^{-x_{i,j}^*})(\pi_{i,j} - \sum_{i \in \mathcal{N}} (P_{i,j}^{U*} \varpi_j + \theta_j Q_i) x_{i,j})$$

$$= \frac{(1 - e^{-x_{i,j}^*}) \sum\limits_{i \in \mathcal{N}} (P_{i,j}^U \varpi_j + \theta_j Q_i)}{\sum\limits_{i \in \mathcal{N}} (P_{i,j}^{U*} \varpi_j + \theta_j Q_i)(1 - e^{-x_{i,j}^*})}$$

$$\times (\pi_{i,j}^f - \sum_{i \in \mathcal{N}} (P_{i,j}^{U*} \varpi_j + \theta_j Q_i)(1 - e^{-x_{i,j}^*}))$$

$$= \pi_{i,j}^f - \sum_{i \in \mathcal{N}} (P_{i,j}^{U*} \varpi_j + \theta_j Q_i)(1 - e^{-x_{i,j}^*})$$

$$= U_j^f, \tag{53}$$

where $U_j^f$ is defined as the fractional utility of mobile user $j$. In the following, we only need to prove that the fractional

utility, i.e., $U_j^f$, is no less than zero. Based on our proposed payment design, $U_j^f$ can be further calculated as

$$U_j^f = \pi_{i,j}^f - (P_{i,j}^{U*} \varpi_j + \theta_j Q_i)(1 - e^{-x_{i,j}^*})$$

$$= \sum_{i \in \mathcal{N}} \sum_{j' \in \mathcal{M}/\{j\}} (a_{i,j'} - P_{i,j'}^{D*} \varpi_{B,j'} - P_{i,j'}^{U*} \varpi_{j'})(1 - e^{-x_{i,j'}^*})$$

$$- \max_{\boldsymbol{X}} \sum_{i \in \mathcal{N}} \sum_{j' \in \mathcal{M}/\{j\}} (a_{i,j'} - P_{i,j'}^D \varpi_{B,j'} - P_{i,j'}^U \varpi_{j'})(1 - e^{-x_{i,j'}})$$

$$- (P_{i,j}^{U*} \varpi_j + \theta_j Q_i)(1 - e^{-x_{i,j}^*})$$

$$= \max_{\boldsymbol{X}} \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} (a_{i,j} - P_{i,j}^{D*} \varpi_{B,j} - P_{i,j}^{U*} \varpi_j)(1 - e^{-x_{i,j}^*})$$

$$\tag{54}$$

$$- \max_{\boldsymbol{X}} \sum_{i \in \mathcal{N}} \sum_{j' \in \mathcal{M}/\{j\}} (a_{i,j'} - P_{i,j'}^D \varpi_{B,j'} - P_{i,j'}^U \varpi_{j'})(1 - e^{-x_{i,j'}})$$

$$\geq 0,$$

where the last inequality holds because the first term in (54) is the maximal social welfare with total $N$ mobile users, while the second term in (54) is the maximal social welfare excluding mobile user $j$. Thus, the fractional payment is individually rational, so is $\mathbb{E}\{U_j(\boldsymbol{L})\}$. This completes the proof. ∎

**Lemma 5.** *The proposed mechanism (OTM) is incentively compatible (truthful) in expectation.*

*Proof.* We assume that mobile user $j$ misreports its actual private information $\varpi_j$ as $\acute{\varpi}_j$, and the fractional solutions based on this false value are $\acute{x}_{i,j}$, $\acute{P}_{i,j}^D$, and $\acute{P}_{i,j}^U$. Then, the fractional utility due to the falsely reported information can be calculated as

$$\acute{U}_j^f = \max_{\boldsymbol{X}} \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} (a_{i,j} - \acute{P}_{i,j}^D \varpi_{B,j} - \acute{P}_{i,j}^U \varpi_j)(1 - e^{-\acute{x}_{i,j}})$$

$$- \max_{\boldsymbol{X}} \sum_{i \in \mathcal{N}} \sum_{j' \in \mathcal{M}/\{j\}} (a_{i,j'} - P_{i,j'}^D \varpi_{B,j'} - P_{i,j'}^U \varpi_{j'})(1 - e^{-x_{i,j'}})$$

The difference between $U_j^f$ and $\acute{U}_j^f$ can be calculated as

$$U_j^f - \acute{U}_j^f = \max_{\boldsymbol{X}} \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} (a_{i,j} - P_{i,j}^{D*} \varpi_{B,j} - P_{i,j}^{U*} \varpi_j)(1 - e^{-x_{i,j}^*})$$

$$- \max_{\boldsymbol{X}} \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} (a_{i,j} - \acute{P}_{i,j}^D \varpi_{B,j} - \acute{P}_{i,j}^U \varpi_j)(1 - e^{-\acute{x}_{i,j}}) \tag{55}$$

$$\geq 0.$$

Since the first term in inequality (55) is the maximal objective under the true values of $\varpi_j$ and $\varpi_{B,j}$, $U_j^f - \acute{U}_j^f \geq 0$ means the fractional payment is truthful. As a result, by combing (53), our proposed mechanism is also truthful in expectation. This completes the proof. ∎

**Theorem 1.** *The proposed mechanism (OTM) is $1 - \frac{1}{e}$-approximation, truthful and individual rationality in expectation.*

*Proof.* The conclusion can be proved by combining the above **Lemmas**. ∎

In the next subsection, we will reconsider the constraint (8), which is the energy budget constraint for each mobile user, and adopt our designed OTM to devise an online truthful mechanism.

## 4.2 Online Truthful Mechanism Design

Since the energy constraint for each mobile user is applied to all time slots, some mobile users may use up their energies before $T$ time slots, so that they cannot participate the competition campaign in the rest time slots. This may result in a lower objective value because it narrows down the possible choosing space of mobile users, and the remaining mobile users may submit valuation, i.e., $c_j^{(t)}$ with high values in the future time slots. To address this issue, we intentionally extend lifetime of mobile users to allow as many mobile users as possible be survived in any time slot, so as to allow the BS to explore more cost-efficient mobile users for increasing the objective value. Following this intuition, we establish a relationship between the remained energy and the total energy, and intentionally decrease the mobile users' winning probabilities in the upcoming time slots if they have been selected before.

Specifically, we introduce two auxiliary variables, i.e., $\nu_j^{(t)}$ and $\nu_{B,j}^{(t)}$ for each mobile user $j$, and increase their values from an initial value $\frac{1}{\varphi}$, where $\varphi = \max\limits_{t \in \mathcal{T}, j \in \mathcal{U}^{sel,(t)}} \frac{e_{i_j,j}^{(t)}}{E_j}$, denotes the maximum ratio between the energy consumption and the total energy. Obviously, $\varphi$ should be much less than 1 as no mobile user is willing to use too much energy within a single time slot. Moreover, instead of applying actual unit power cost $c_j^{(t)}$ and $c_B^{(t)}$ in the **Algorithm** 2, we replace them by $\varpi_j^{(t)} = \beta_j c_j \varphi \nu_j^{(t-1)}$ and $\varpi_{B,j}^{(t)} = \beta_j c_B \varphi \nu_{B,j}^{(t-1)}$, where $\beta_j$ $(\beta_j > 1)$, called the inflating factor, is used to improve the performance of our online truthful mechanism, and $\beta = \max\limits_{j \in \mathcal{M}} \beta_j$. Note that $\nu_j^{(t-1)}$ and $\nu_{B,j}^{(t-1)}$ are both updated carefully to ensure that the mobile user who has lower remaining energy will have a less chance to be selected in future time slots. For convenience, we summarize the online truthful mechanism as in the **Algorithm** 3. It is worth noting that by using the updating manner on $\nu_j^{(t-1)}$ and $\nu_{B,j}^{(t-1)}$ as in **Algorithm** 3, theoretical CR exists. Moreover, unlike [30], [51], where the formulated problems were linear, in this paper, we consider a nonlinear scenario, but still having a CR guarantee.

**Lemma 6.** *The CR of the proposed online truthful mechanism is $\beta(1 - \frac{1}{e})(2^\varphi - 1)$ in expectation.*

*Proof.* Let's begin with the proof of the following inequalities

$$\nu_j^{(t)} \approx \frac{1}{\varphi}(2^{\frac{\sum_{t \in \mathcal{T}} e_{i_j,j}^{(t)}}{E_j}} - 1) \le \frac{1}{\varphi}(2^\varphi - 1), \qquad (56)$$

$$\nu_{B,j}^{(t)} \approx \frac{1}{\varphi}(2^{\frac{\sum_{t \in \mathcal{T}} e_{i_j,j}^{(t)}}{E_j}} - 1) \le \frac{1}{\varphi}(2^\varphi - 1). \qquad (57)$$

We prove inequalities (56) and (57) by induction on time index, $t$. Initially, we assume those inequalities are all held at time slot $t-1$, and the task $i$ will be assigned to the mobile user $j$ at time slot $t$. By taking inequality (56) as an example,

---

**Algorithm 3:** Online truthful mechanism.

1 **Initialization**;
2 $\nu_j^{(0)} = \nu_{B,j}^{(0)} = \frac{1}{\varphi}, \; \forall j \in \mathcal{M}$;
3 **for** *each time slot* $t \in \mathcal{T}$ **do**
4      $\varpi_j^{(t)} = \beta c_j^{(t)} \varphi \nu_j^{(t-1)}$;
5      $\varpi_{B,j}^{(t)} = \beta c_B^{(t)} \varphi \nu_{B,j}^{(t-1)}$;
6      Run **Algorithm** 2 to get the set of selected mobile users and their indexes of execution tasks, which are denoted as $\mathcal{U}^{sel,(t)}$ and $i_j$, respectively;
7      Based on (12), calculate energy consumptions $e_{i_j,j}^{(t)}$ at mobile user $j \in \mathcal{U}^{sel,(t)}$;
8      $\nu_j^{(t)} = \nu_j^{(t-1)}(1 + \frac{e_{i_j,j}^{(t)}}{E_j}) + \frac{e_{i_j,j}^{(t)}}{\varphi E_j}, \; j \in \mathcal{U}^{sel,(t)}$;
9      $\nu_{B,j}^{(t)} = \nu_{B,j}^{(t-1)}(1 + \frac{e_{i_j,j}^{(t)}}{E_j}) + \frac{e_{i_j,j}^{(t)}}{\varphi E_j}, \; j \in \mathcal{U}^{sel,(t)}$;
10      $\nu_j^{(t)} = \nu_j^{(t-1)}, \nu_{B,j}^{(t)} = \nu_{B,j}^{(t-1)}, \; j \in \mathcal{M}/\mathcal{U}^{sel,(t)}$;

---

we have

$$\nu_j^{(t)} = \nu_j^{(t-1)}(1 + \frac{e_{i_j,j}^{(t)}}{E_j}) + \frac{e_{i_j,j}^{(t)}}{\varphi E_j},$$

$$\approx \frac{1}{\varphi}(2^{\frac{\sum_{t \in \mathcal{T}/t} e_{i,j}^{(t)}}{E_j}} - 1)(1 + \frac{e_{i_j,j}^{(t)}}{E_j}) + \frac{e_{i_j,j}^{(t)}}{\varphi E_j},$$

$$= \frac{1}{\varphi}(2^{\frac{\sum_{t \in \mathcal{T}/t} e_{i,j}^{(t)}}{E_j}}(1 + \frac{e_{i_j,j}^{(t)}}{E_j}) - 1),$$

$$\approx \frac{1}{\varphi}(2^{\frac{\sum_{t \in \mathcal{T}} e_{i_j,j}^{(t)}}{E_j}} - 1) \le \frac{1}{\varphi}(2^\varphi - 1), \qquad (58)$$

where the approximation in (58) holds because $2^x \approx 1 + x$ when $x$ is small. Let $\overline{P}^{(t)}$ be the optimal objective value of the problem $[\mathcal{RP}]$ at time slot $t$ without the consideration of the constraint (8), and $P^{(t)*}$ is the optimal objective value at time slot $t$ calculated by the **Algorithm** 2. We have

$$P^{(t)*} = \sum_{j \in \mathcal{U}^*} (a_{i_j,j} - P_{i_j,j}^{D*} \varpi_{B,j} - P_{i_j,j}^{U*} \varpi_j)$$

$$\ge \sum_{j \in \overline{\mathcal{U}}} (a_{i_j,j} - \overline{P}_{i_j,j}^D \varpi_{B,j} - \overline{P}_{i_j,j}^U \varpi_j)$$

$$\ge \beta(2^\varphi - 1) \sum_{j \in \overline{\mathcal{U}}} (a_{i_j,j} - \overline{P}_{i_j,j}^D c_B - \overline{P}_{i_j,j}^U c_j)$$

$$= \beta(2^\varphi - 1)\overline{P}^{(t)}. \qquad (59)$$

Therefore, the objective value $\hat{P}^{(t)}$ from the **Algorithm** 3 at time slot $t$ can be calculated as

$$\hat{P}^{(t)} = \sum_{j \in \mathcal{U}^{sel,(t)}} (a_{i_j,j} - P_{i_j,j}^D c_B - P_{i_j,j}^U c_j)$$

$$\ge \sum_{j \in \mathcal{U}^{sel,(t)}} (a_{i_j,j} - P_{i_j,j}^D \varpi_{B,j} - P_{i_j,j}^U \varpi_j) \qquad (60)$$

$$\ge (1 - \frac{1}{e})P^{(t)*}$$

$$\ge \beta(1 - \frac{1}{e})(2^\varphi - 1)\overline{P}^{(t)},$$

where inequality (60) holds since $c_j \leq \varpi_j$ and $c_B \leq \varpi_{B,j}$. Summing over total time slots, we have

$$\hat{P} = \sum_{t \in \mathcal{T}} \hat{P}^{(t)} \geq \beta(1 - \frac{1}{e})(2^\varphi - 1) \sum_{t \in \mathcal{T}} \overline{P}^{(t)}$$

$$\geq \beta(1 - \frac{1}{e})(2^\varphi - 1)P^{opt}, \quad (61)$$

where $P^{opt}$ is the optimal objective value of the problem $[\mathcal{RP}]$, and the last inequality holds because $\sum_{t \in \mathcal{T}} \overline{P}^{(t)}$ is the summation of all independent time slots without the consideration of the constraint (8), which is no less than that in our proposed scenario. This completes the proof. ∎

**Lemma 7.** *Our proposed online truthful mechanism can obtain the result in polynomial time.*

*Proof.* In this paper, the computational complexity of the proposed online mechanism is evaluated in terms of computation times with respect to the number of offloading tasks, mobile users, and total time slots. Remind that our proposed online truthful mechanism consists of two parts, Algorithm 2 and Algorithm 3. As for Algorithm 2, we solve the fractional convex problem $[\mathcal{EPRO}1]$ by gradient descent optimizing solver. Suppose that $L_{max}$ and $I_{max}$ are the maximal iteration numbers to solve the Lagrangian primal problem and its primal problem, respectively. Then, the computational complexity of Algorithm 2 can be calculated as $O(NM + L_{max}I_{max})$. For the Algorithm 3, the total computational complexity is $O(|\mathcal{T}| \times (NM + L_{max}I_{max}))$. In summary, the proposed online truthful mechanism runs in polynomial time, which completes the proof. ∎

**Theorem 2.** *The proposed online truthful mechanism is truthful and individually rational, and the CR is $\beta(1 - \frac{1}{e})(2^\varphi - 1)$ in expectation.*

*Proof.* The conclusion can be proved by combining the **Lemma** 6 and the **Theorem** 1. ∎

## 5 NUMERICAL RESULTS

In this section, numerical simulations are conducted to verify the effectiveness of our proposed online truthful mechanism. Since the total social welfare, and the utility of mobile users are the most important economical metrics and the competitive ratio is vital to measure an online truthful mechanism, we will focus on evaluating these three performance metrics with respect to different numbers of offloading tasks and mobile users. In the simulations, the wireless channels between the BS and mobile users experience Rayleigh fading and all channel coefficients are zero-mean, circularly symmetric complex Gaussian (CSCG) random variables with variances $d^{-\frac{v}{2}}$, where $d$ is the distance between the transmitter and the receiver and $v = 4$. Table 3 lists the main simulation parameter values, most of which have been employed in [6], [25], [27], [38], [53]. In the following figures, each performance point is derived by averaging 300 independent runs. For comparison purpose, the following three benchmarks online strategies are simulated as well.

- Lyapunov based online mechanism (LOM): The Lyapunov optimization was used in [27] for designing

### Table 3
### MAIN SIMULATION PARAMETERS

| Parameter | Value |
|---|---|
| Cell radius | 500 m |
| Download link bandwidth | 40 MHz |
| Upload link bandwidth | 10 MHz |
| $P_B^{max}$ at the BS | 46 dBm |
| $P_j^{max}$ at mobile users | 23 dBm |
| Background noise average power | -60 dBm |
| Total running time | 20 minutes |
| Time slot length | 60 seconds |
| Input task size | Randomly from 10 to 30 MB |
| Output task size | 20% of the input data |
| CPU cycles coefficient | Randomly between [125, 375] cycles/Byte |
| $\beta_j$ at mobile users | $10^{-26}$ |
| $\theta_B$ at the BS | $\$10^{-10}$ |
| $c_B^{(t)}$ at the BS | $\$0.1$ |
| Benefits, i.e., $r_{i,j}^{(t)}$, at the BS | Randomly over (1, 2] |
| $\theta_j$ at mobile users | $\$0.5 \times 10^{-10}$ |
| $c_j^{(t)}$ at mobile users | Randomly from [\$0.5, \$1] |
| $D_i^{(t)}$ delay demand | Randomly from [1, 15] seconds |
| $F_j$ at mobile users | 2 GHz |

online mechanism, but working at a fixed maximal upload transmission power. Note that since the objective by using the Lyapunov method should be a long-term one, we only calculate the first $|\mathcal{T}|$ time slots in the simulation.

- Fixed power online mechanism (FPOM): The transmission power of both upload and download links are fixed in advance, and the online mechanism problem is solved by our proposed method.

- Random online mechanism (ROM): In each time slot, the optimization variables, i.e., $\boldsymbol{X}$, are determined randomly, and mobile users will not be excluded from future participation until their energy budgets deplete.
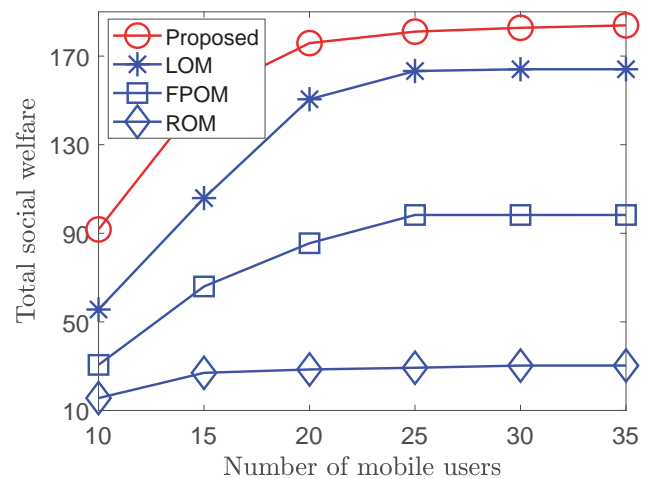


Figure 2. Total social welfare versus number of mobile users when offloaded tasks are with different numbers in each time slot.

Fig. 2 illustrates the total social welfare obtained by different online mechanisms with respect to different number of mobile users when the tasks offloaded by the BS are randomly chosen from [10, 20] for each time slot, and the
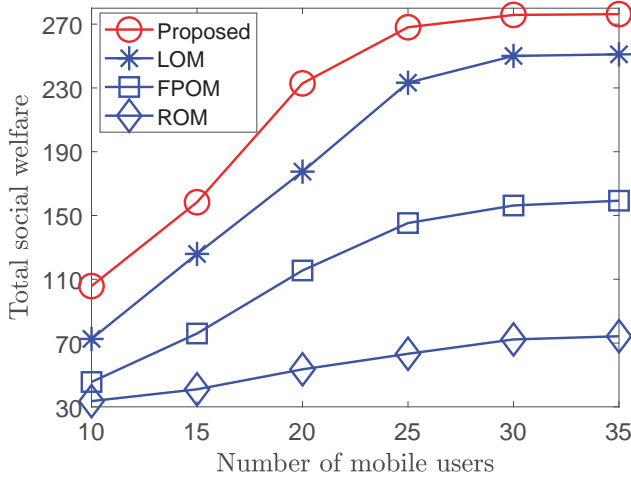
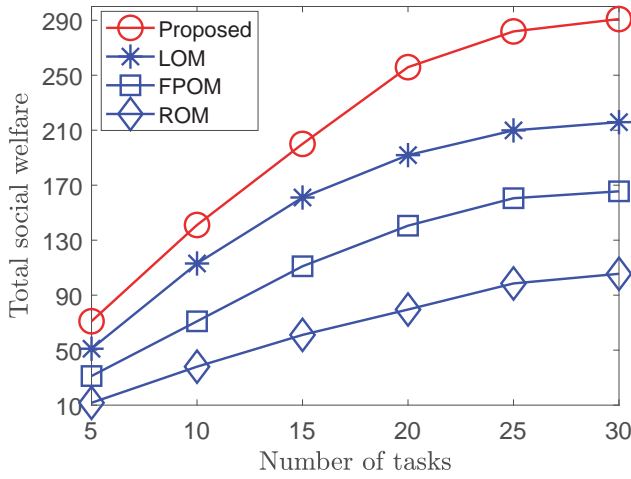Figure 3. Total social welfare versus number of mobile users.



Figure 5. Total social welfare versus maximal transmission power at mobile users.



Figure 4. Total social welfare versus number of offloaded tasks.



Figure 6. The profits versus number of mobile users.

energy budget is set to $8 \times 10^3$ (Joule). It can be seen from this figure that the total social welfare by all online mechanisms increases with the number of mobile users till reaching saturation when the number of mobile users is large enough. This can be explained as follows. With the number of mobile users increasing, more and more tasks can be successfully offloaded to mobile users so that the total social welfare increases. However, with the excessive amount of mobile users, e.g., 30, all tasks offloaded by the BS are accepted and executed by winning mobile users, and no extra mobile users can contribute to the total social welfare. Moreover, our proposed online mechanism is superior to all other three online mechanisms. This is because our proposed method jointly optimizes upload and download links' transmission powers, while the LOM ignores the optimization of upload link transmission power, and the FPOM doesn't consider the power control. Furthermore, Fig. 3 reevaluates the relationship between the total social welfare and a various number of mobile users when there are 18 tasks needing to be offloaded in each time slot. From this figure, we can almost draw the same conclusions as those in Fig. 3. In addition, it is worth noting that when the number of mobile users is large enough, the performance by the LOM is close
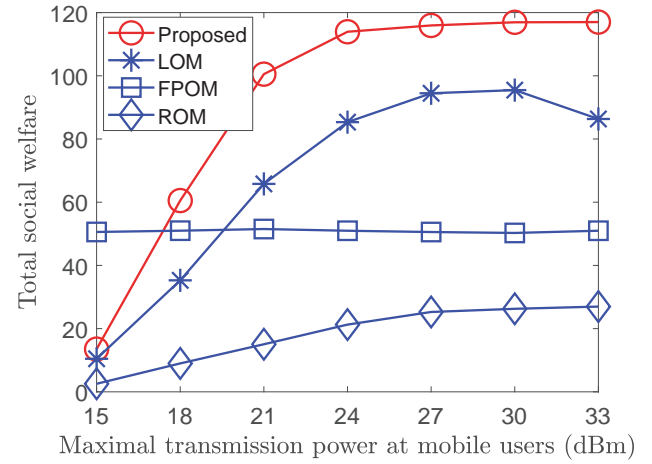
to our proposed one. The reason behind could be that more mobile users mean that more potential suitable mobile users (e.g., larger $r_{i,j}$ and lower bidding cost $c_{i,j}^{(t)}$) can be selected by the BS, which, to some extent, can reduce the effect of upload link power control.

Fig. 4 reveals the total social welfare under different numbers of arriving tasks in each time slot. In the simulation, the total number of mobile users is 18, i.e., $|\mathcal{M}| = 18$, and the energy budget is set to $8 \times 10^3$ (Joule). From this figure, we can see that the achievable total social welfare will reach a plateau at the end. This is because with the number of tasks increasing, each task is always served by its selected mobile user while other unselected tasks cannot increase the achievable total social welfare. Note that when the number of tasks is large enough, the total social welfare is still gradually increasing. It can be explained as follows. Since newly arrived tasks have different deadline requirements, the BS will prefer less delay intensive tasks so as to lower its communication and computation resource demand, resulting in the slow increase of the total social welfare. Moreover, this figure also verifies the superiority of our proposed online mechanism over others, which further demonstrates
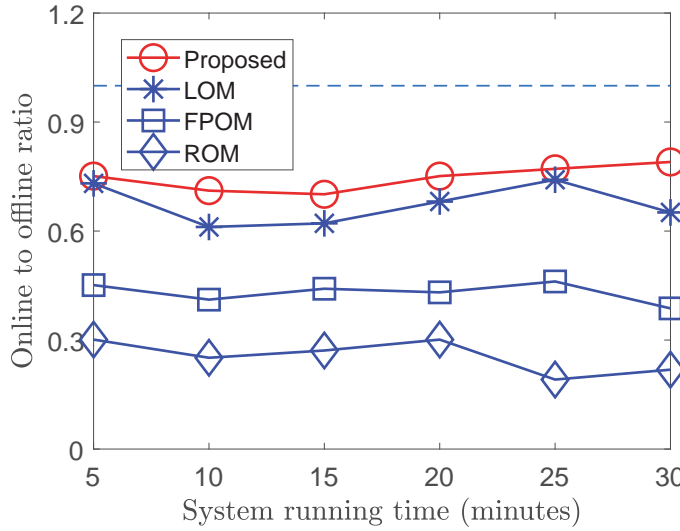
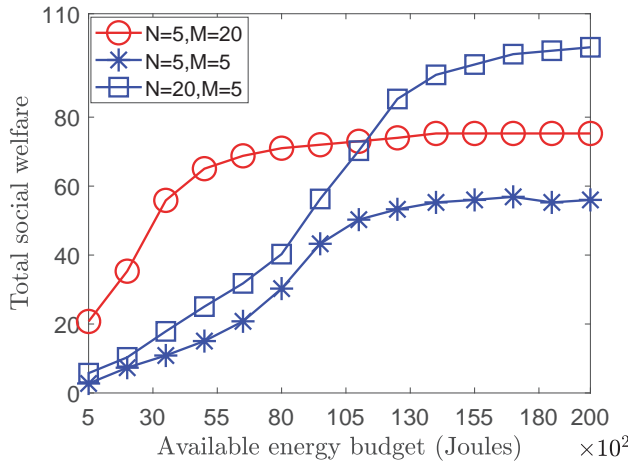Figure 7. Online to offline ratio versus system running time.



Figure 9. Total social welfare versus unit power cost at the BS.



Figure 8. Total social welfare versus available energy budget at mobile users.



Figure 10. Algorithm execution time for one time slot and total time slots.

that our proposed mechanism can always achieve better performance under any combination of offloaded tasks and mobile users.

Fig. 5 depicts the relationship between the total social welfare and the maximal transmission power at mobile users. It can be observed that the total social welfare by the proposed online mechanism increases first and then keeps stable when the maximal transmission power at mobile users is large enough. This is because, in order to meet the delay requirements of tasks, mobile users have to transmit back the computation results with their maximal but very limited transmission powers. With the increase of the maximal transmission power, more tasks' delay requirements can be satisfied, which leads to an increase in the total social welfare. However, when the maximal transmission power is large enough, e.g., $P_j^{max} = 35$ dBm, the optimal transmission power becomes less than the maximal one. In addition, since all tasks have been successfully offloaded, the total social welfare becomes almost a constant. Furthermore, the total social welfare by the LOM method first increases to
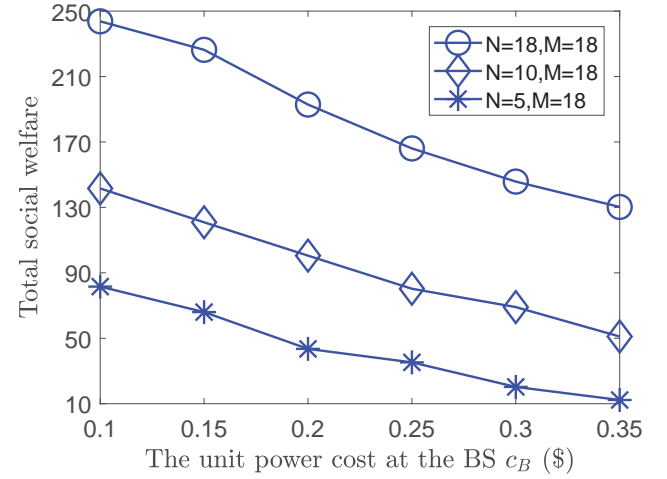
a maximum, and then reduces slightly when the maximal transmission power is over 35 dBm. This decreasing trend is intuitive. Since the LOM always works at a maximal upload link transmission power, according to the definition of our objective function, the total social welfare decreases when the total social welfare reaches its maximum at almost 32 dBm. In addition, the total social welfare achieved by the FPOM stays almost unchanged. This is because for the FPROM method, both its upload and download links' transmission powers are fixed so that the maximal transmission power at mobile users has no effects on the total social welfare.

Fig. 6 shows the trend of average profit, i.e., average utilities of mobile users and the BS, with the number of mobile users under different online mechanisms. From this figure, we can observe that both utilities of mobile users and the BS increase with the number of mobile users. This is because more mobile users can accept more offloaded tasks for execution, which increases the benefits. Moreover, since our proposed online mechanism jointly considers the communication and computation resource allocations, it is obvious that our proposed scheme can gain more benefits

compared to others. Furthermore, the average utility of mobile users by our proposed online mechanism is no less than zero, which manifests the property of individual rationality holds.

Fig. 7 presents the comparisons among different online mechanisms with respect to online to offline ratio (i.e., CR ratio) along system running time. Note that the optimal offline solution is obtained by the brute force method, and to reduce the computation time, we limit both the numbers of mobile users and tasks to be 5. Also, note that by letting $\varphi = 0.1$ and $\beta = 10$, we have the theoretical CR about 0.45, which is the worst-case performance. However, from this figure, we can observe that the actually achievable CR by our proposed mechanism is around 0.75, which is much better than this worst-case CR. Moreover, the actual CR almost stays unchanged along system running time, which shows that the proposed online mechanism is robust to the running time.

Fig. 8 reveals the relationship between the total social welfare and the available energy budget at mobile users with different number of mobile user and task pairs. Three cases are under consideration. In case one, the number of tasks is 5 and the number of mobile users is 20; in case two, the number of tasks is 20 and the number of mobile users is 5; and in case three, the numbers of both mobile users and tasks are 5. From this figure, it can be seen that the total social welfare first increases and eventually strikes a balance regardless of the amounts of mobile users and tasks. This is because, with the increase of the available energy budget, more mobile users can survive to execute more tasks so that the total social welfare increases. However, since there are only at most five tasks for executions, the total social welfare will not continue increasing when the available budget is large enough. Besides, for case one, the total social welfare experiences a surge in the early increasing stage with the energy budget. This is because when the energy budget increases a little bit, twenty mobile users could be enough to alternatively execute five offloaded tasks. Moreover, the total social welfare under case two, is always larger than that under case three. The reason behind this can be explained as follows. In case two, there is a sufficient number of tasks so that there will always be tasks with less stringent delay tolerance requirements to choose from. Hence, both mobile users and the BS could consume less computation and communication resources for offloaded task execution so that the total social welfare becomes higher. In addition, the total social welfare by case two surpasses that by case one when the available energy budget becomes large enough. This is because compared to the diversity of mobile users, the diversity of tasks will contribute more to the increase of the total social welfare, which can also be corroborated by comparing Fig. 3 and Fig. 4 in the stable state.

Fig. 9 shows the effect of unit power cost at the BS on total social welfare. From this figure, it can be seen that the total social welfare decreases with the increase of unit power cost at the base station. This is because the base station will try to select the tasks with smaller data sizes and mobile users with good channels to offload, as the larger data size tasks and bad channels could cause the increase in transmission power at the base station in order to meet the deadline requirement. As a result, less tasks and mobile

users can be selected for offloading, which in turn causes the decrease in total social welfare. In addition, Fig. 10 presents the execution time of our proposed algorithm in one time slot and total time slots, respectively. Obviously, the time consumption in each time slot for algorithm execution increases with the number of offloaded tasks, but is far less than 5 seconds. Since the BS in practice has much more powerful computation capacity, the actual execution time could be much less than ours.

## 6 CONCLUSION

In this paper, a nonlinear online truthful mechanism for task offloading in edge computing systems has been proposed. By considering the facts that the arrival of tasks is dynamic and each mobile user is energy-constrained in practice, we formulate a social welfare maximization problem by jointly considering task offloading decisions, and both computation and communication resource allocation. We convert this time coupled incentive mechanism design problem into several one-shot ones, and solve them by our proposed IRSM framework. Finally, we rationally combine the results of one-shot incentive mechanism to design a nonlinear online incentive mechanism. Theoretical analyses guarantee the properties of IR, IC, and computational efficiency. Numerical results further demonstrate the effectiveness and superiority of our proposed nonlinear online truthful mechanism.

## REFERENCES

[1] X. Li, R. Lu, X. Liang, X. Shen, J. Chen, and X. Lin, "Smart community: an internet of things application," *IEEE Commun. Mag.*, vol. 49, no. 11, pp. 68-75, Nov. 2011.

[2] (2017) Hyrax crowd-sourcing mobile devices to develop edge clouds. [Online]. Available: http://hyrax.dcc.fc.up.pt.

[3] N. Nisan, T. Roughgarden, E. Tardos, and V. V. Vazirani, "Algorithmic Game Theory," *Cambridge University Press*, 2007.

[4] Y. Zhao, S. Zhou, T. Zhao, and Z. Niu. "Energy-efficient task offloading for multiuser mobile cloud computing," *in Proc. IEEE/CIC Int. Conf. Commun. China (ICCC)*, Shenzhen, China, 2015, pp. 1-5.

[5] L. Yang, B, Liu, J. Cao, etal. "Joint Computation Partitioning and Resource Allocation for Latency Sensitive Applications in Mobile Edge Clouds," *in Proc. IEEE Int. Conf. Cloud. Computing*, Honolulu, CA, USA, 2017, pp. 246-254.

[6] T. Thinh, J. Tang, Q. La and T. Quek, "Offloading in Mobile Edge Computing: Task Allocation and Computational Frequency Scaling," *IEEE Transactions on Communications*, pp. 1-1, 2017.

[7] S. Sardellitti, S. Barbarossa, and G. Scutari, "Distributed mobile cloud computing: Joint optimization of radio and computational resources," *in Proc. IEEE Globecom Workshops (GC Wkshps)*, Austin, TX, USA, 2014, pp. 1505-1510.

[8] M. H. Chen, M. Dong, and B. Liang, "Joint offloading decision and resource allocation for mobile cloud with computing access point," *in Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, Pudong, China, 2016, pp. 3516-3520.

[9] J. Du, L. Zhao, J. Feng and X. Chu, "Computation Offloading and Resource Allocation in Mixed Fog/Cloud Computing Systems With Min-Max Fairness Guarantee," *IEEE Transactions on Communications*, vol. 66, no. 4, pp. 1594-1608, April 2018.

[10] G. Zhang, W. Zhang, Y. Cao, D. Li and L. Wang, "Energy-Delay Tradeoff for Dynamic Offloading in Mobile-Edge Computing System With Energy Harvesting Devices," *IEEE Transactions on Industrial Informatics*, vol. 14, no. 10, pp. 4642-4655, Oct. 2018.

[11] M. Kamoun, W. Labidi, and M. Sarkiss, "Joint resource allocation and offloading strategies in cloud enabled cellular networks," *in Proc. IEEE Int. Conf. Commun. (ICC)*, London, U.K., 2015, pp. 5529-5534.

[12] W. Labidi, M. Sarkiss, and M. Kamoun, "Joint multi-user resource scheduling and computation offloading in small cell networks," *in Proc. IEEE Int. Conf. Wireless Mobile Comput. Netw. Commun. (WiMob)*, Abu Dhabi, UAE, 2015, pp. 794-801.

[13] Y. Mao, J. Zhang, and K. B. Letaief, "Dynamic computation offloading for mobile-edge computing with energy harvesting devices," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 12, pp. 3590-3605, Dec. 2016.

[14] Y. Mao, J. Zhang, and K. B. Letaief, "Stochastic Joint Radio and Computational Resource Management for Multi-User Mobile-Edge Computing Systems," *IEEE Wirles. Commun.*, vol. 16, no. 9, pp. 5994-6009, Sep. 2017.

[15] X. Qiu, L. Liu, W. Chen, Z. Hong, and Z. Zheng, "Online Deep Reinforcement Learning for Computation Offloading in Blockchain-Empowered Mobile Edge Computing," *IEEE Trans. veh. Technol.*, vol. 68, no. 8, pp. 8050-8062, Aug. 2019.

[16] L. Huang, S. Bi, and Y. A. Zhang, "Deep Reinforcement Learning for Online Computation Offloading in Wireless Powered Mobile-Edge Computing Networks," *IEEE Trans. Mobile Comput.*, vol. 19, no. 11, pp. 2581-2593, Nov. 2020.

[17] Z. Zhou, Q. Wu and X. Chen, "Online Orchestration of Cross-Edge Service Function Chaining for Cost-Efficient Edge Computing," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 8, pp. 1866-1880, Aug. 2019.

[18] A. Bozorgchenani, S. Maghsudi, D. Tarchi and E. Hossain, "Computation Offloading in Heterogeneous Vehicular Edge Networks: On-line and Off-policy Bandit Solutions," *IEEE Transactions on Mobile Computing*, 2021.

[19] A. Kiani, and N. Ansari, "Towards Hierarchical Mobile Edge Computing: An Auction-Based Profit Maximization Approach," *IEEE Internet of Things Journal*, vol. PP, no. 99, pp. 1-10, Sep. 2017.

[20] L. U. Khan et al., "Federated Learning for Edge Networks: Resource Optimization and Incentive Mechanism," *IEEE Commun. Mag.*, vol. 58, no. 10, pp. 88-93, Oct. 2020.

[21] M. LiWang, et al., "A Truthful Reverse-Auction Mechanism for Computation Offloading in Cloud-enabled Vehicular Network," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 4214-4227, Jun. 2019.

[22] Z. Chang, W. Guo, X. Guo, Z. Zhou, and T. Ristaniemi, "Incentive Mechanism for Edge-Computing-Based Blockchain," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 11, pp. 7105-7114, Nov. 2020.

[23] F. Mashhadi, S. A. S. Monroy, and A. Bozorgchenani, et al, "Optimal auction for delay and energy constrained task offloading in mobile edge computing," *Comput. Networks*, vol. 183, pp. 1-10, 2020.

[24] Q. Wang, S. Guo, J. Liu, C. Pan and L. Yang, 'Profit Maximization Incentive Mechanism for Resource Providers in Mobile Edge Computing,' *IEEE Transactions on Services Computing*, to appear.

[25] G. Li and J. Cai, "An Online Incentive Mechanism for Collaborative Task Offloading in Mobile Edge Computing," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 624-636, Jan. 2020.

[26] J. He, D. Zhang, Y. Zhou and Y. Zhang,"A Truthful Online Mechanism for Collaborative Computation Offloading in Mobile Edge Computing," *IEEE Transactions on Industrial Informatics*, vol. 16, no. 7, pp. 4832-4841, Jul. 2020.

[27] D. Zhang et al., "Near-optimal and Truthful Online Auction for Computation Offloading in Green Edge-Computing Systems," *IEEE Trans. Mobile Comput.*, vol. 19, no. 4, pp. 880-893, Apr. 2020.

[28] G. Li, J. Cai and H. Chen, "Online Truthful Mechanism Design in Wireless Communication Networks," *IEEE Wireless Communications*, vol. 28, no. 4, pp. 159-165, August 2021.

[29] S. Dobzinski and S. Dughmi, "On the Power of Randomization in Algorithmic Mechanism Design," *in Proc. of IEEE FOCS*, 2009.

[30] G. Li and J. Cai,"An Online Mechanism for Crowdsensing with Uncertain Task Arriving," *in Proc. IEEE ICC*, Kansas, MO, 2018, pp. 1-6.

[31] G. Li and J. Cai, "An Online Incentive Mechanism for Crowdsensing with Random Task Arrivals," *IEEE Internet Things J.*, vol. 7, no. 4, pp. 1-14, Jan. 2020.

[32] D. Zhao, X. Y. Li, and H. D. Ma, "Budget-feasible online incentive mechanisms for crowdsourcing tasks truthfully," *IEEE Trans. Net.*, vol. 24, no. 2, pp. 647-661, Apr. 2016.

[33] P. A. Apostolopoulos, G. Fragkos, E. E. Tsiropoulou and S. Papavassiliou, "Data Offloading in UAV-assisted Multi-access Edge Computing Systems under Resource Uncertainty," *IEEE Transactions on Mobile Computing*, 2021.

[34] C. Yi, S. Huang, and J. Cai, "An incentive mechanism integrating joint power, channel and link management for social-aware D2D content sharing and proactive caching," *IEEE Trans. Mobile Comput.*, vol. 17, no. 4, pp. 789-802, Apr. 2018.

[35] P. Wang, Z. Zheng, B. Di and L. Song, "HetMEC: Latency-Optimal Task Assignment and Resource Allocation for Heterogeneous Multi-Layer Mobile Edge Computing," *IEEE Trans. Wireless Commun.*, vol. 18, no. 10, pp. 4942-4956, Oct. 2019.

[36] J. Zhang, X. Huang and R. Yu, "Optimal Task Assignment With Delay Constraint for Parked Vehicle Assisted Edge Computing: A Stackelberg Game Approach," *IEEE Commun. Lett.*, vol. 24, no. 3, pp. 598-602, March 2020.

[37] Y. Zhang, L. Song, W. Saad, Z. Dawy, and Z. Han, "Contract based incentive mechanisms for device-to-device communications in cellular networks," *IEEE J. Sel. Areas Commun.*, vol. 33, no. 10, pp. 2144-2155, Oct. 2015.

[38] L. Pu, X. Chen, J. Xu, and X. Fu, "D2D fogging: An energy-efficient and incentive-aware task offloading framework via network assisted collaboration," *IEEE J. Sel. Areas Commun.*, vol. 34, no. 12, pp. 3887-3901, Dec. 2016.

[39] 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, Aug. 2017.

[40] A. P. Miettinen and J. K. Nurminen, "Energy efficiency of mobile clients in cloud computing,"*in Proc. 2nd USENIX Conf. Hot Topics Cloud Comput. (HotCloud)*, 2010, pp. 1-4.

[41] S. Melendez and M. P. McGarry, "Computation offloading decisions for reducing completion time,"*in Proc. 14th IEEE Annu. Consum. Commun. Netw. Conf. (CCNC)*, Jan. 2017, pp. 160-164.

[42] Y. Wen, W. Zhang, and H. Luo, "Energy-optimal mobile application execution: Taming resource-poor mobile devices with cloud clones," *in Proc. IEEE INFOCOM*, 2012, pp. 2716-2720.

[43] A. Bulut, and T. K. Ralphs, "On the Complexity of Inverse Mixed Integer Linear Optimization," *ArXiv*, vol. abs/2104.09002, 2021.

[44] C.H. Papadimitriou, and M. Yannakakis "The complexity of facets (and some facets of complexity)," *Journal of Computer and System Sciences*, vol. 28, no. 2, pp. 244-259, 1984.

[45] R. Rajaraman, "Randomized Rounding," *In: Kao MY. (eds) Encyclopedia of Algorithms*. Springer, Boston, MA, 2008.

[46] Dughmi, Shaddin, et al, "Optimal Mechanisms for Combinatorial Auctions and Combinatorial Public Projects via Convex Rounding,"*Journal of the ACM*, vol. 63, no. 4, pp. 1-33, Sep. 2016.

[47] N. Nisan, and A. Ronen, "Computationally Feasible VCG Mechanisms," *J.Artif. Intell. Res.(JAIR)* vol. 29, pp. 19-47, 2007.

[48] X. Zhang, C. Wu, Z. Li, and F. C. M. Lau, "A Truthful $(1 - \epsilon)$-Optimal Mechanism for On-demand Cloud Resource Provisioning," *IEEE Transactions on Cloud Computing*, vol. 8, no. 3, pp. 735-748, Sept. 2020.

[49] R. Zhou, Z. Li, C. Wu and M. Chen, "Demand Response in Smart Grids: A Randomized Auction Approach," *IEEE J. Sel. Areas Commun.*, vol. 33, no. 12, pp. 2540-2553, Dec. 2015.

[50] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.

[51] R. Zhou, Z. Li and C. Wu, "A Truthful Online Mechanism for Location-Aware Tasks in Mobile Crowd Sensing," *IEEE Trans. Mobile Comput.*, vol. 17, no. 8, pp. 1737-1749, Aug. 2018.

[52] S. Dughmi, and T. Roughgarden, "Black-box randomized reductions in algorithmic mechanism design," *In Proc. of the 51st IEEE Symposium on Foundations of Computer Science*, 2010, 775-784.

[53] H. Yu, H. D. Tuan, T. Q. Duong, H. V. Poor and Y. Fang, "Optimization for Signal Transmission and Reception in a Macrocell of Heterogeneous Uplinks and Downlinks,"*IEEE Trans. Commun.*, vol. 68, no. 11, pp. 7054-7067, Nov. 2020.