

Truthful Computation Offloading Mechanisms for Edge Computing

Weibin Ma

Department of Computer and
Information Sciences
University of Delaware
Newark, Delaware, USA
Email: weibinma@udel.edu

Lena Mashayekhy

Department of Computer and
Information Sciences
University of Delaware
Newark, Delaware, USA
Email: mlena@udel.edu

Abstract—Edge computing (EC) is a promising paradigm providing a distributed computing solution for users at the edge of the network. Preserving satisfactory quality of experience (QoE) for users when offloading their computation to EC is a non-trivial problem. Computation offloading in EC requires jointly optimizing access points (APs) allocation and edge service placement for users, which is computationally intractable due to its combinatorial nature. Moreover, users are self-interested, and they can misreport their preferences leading to inefficient resource allocation and network congestion. In this paper, we tackle this problem and design a novel mechanism based on algorithmic mechanism design to implement a system equilibrium. Our mechanism assigns a proper pair of AP and edge server along with a service price for each new joining user maximizing the instant social surplus while satisfying all users' preferences in the EC system. Declaring true preferences is a weakly dominant strategy for the users. The experimental results show that our mechanism outperforms user equilibrium and random selection strategies in terms of the experienced end-to-end latency.

Index Terms—Edge Computing, Access Point Allocation, Service Placement, Pricing, Algorithmic Mechanism Design.

I. INTRODUCTION

With the explosive growth of smart devices, a bulk of computationally intensive applications, as exemplified by face recognition, online gaming, and video streaming, are becoming prevalent. However, the smart devices possess limited resources (e.g., limited computation capabilities and battery lifetime), which may lead to unsatisfactory computation experience. The commonly used approach is to offload computational tasks to a powerful cloud platform [1]. However, the long distance between devices and the cloud will cause a significant increase in delay and network congestion.

To overcome this challenge, edge computing (EC) has recently been introduced as an emerging solution that enables offloading computational tasks to the physically proximal EC mini-datacenters, called cloudlets [2], [3]. While EC brings many opportunities to guarantee quality of experience (QoE) for users, new challenges arise due to the restricted coverage of cloudlets and their limited computational resources. To maintain the QoE of users, designing efficient realtime computation offloading is hence becoming crucial in edge computing. The computation offloading problem consists of jointly optimizing access points (APs) allocation and edge service placement

for EC users, which is computationally intractable due to its combinatorial nature.

In this paper, we design a novel mechanism called computation offloading and pricing mechanism (COPM) to satisfy QoE of each joining user by meeting its application-specific end-to-end latency requirements. The goal of our proposed mechanism is to maximize the instant social surplus, which is defined as the sum of the valuation of the new user and the system. To tackle the complexity of COPM, we then propose an online offloading mechanism, called DAPA. When a new user requests an edge service at any time, DAPA collects current information of the system and then assigns an optimal decision pair (the best AP for connection and the best edge server for computation) to the user. It also determines the user's corresponding payment for the edge service. If no feasible solution exists for this user, DAPA suggests the new user offloading its task to the remote cloud.

When a new user requests an edge service, it will report its maximum tolerable end-to-end latency in order to receive the best decision pair to offload and complete its task. A user may misreport this value to increase its own utility. Such an action could inversely decrease the overall system efficiency. Therefore, designing an *incentive-compatible* (or truthful) mechanism in which users have no incentive to lie about their true preferences is extremely important for achieving system efficiency and implementing a system equilibrium. Our goal is to design an efficient incentive-compatible mechanism to determine an optimal decision pair with a corresponding payment for each user satisfying their QoE requirements while maximizing the instant social surplus. To the best of our knowledge, this is the first work that simultaneously optimizes online AP allocation, service placement, and pricing of computation offloading by utilizing algorithmic mechanism design. Our proposed mechanism implements a weakly dominant strategy equilibrium for users.

The rest of the paper is organized as follows. Section II reviews related work. The system model is described in Section III. The problem formulation and COPM mechanism are presented in Section IV. In Section V, we describe our efficient online algorithmic solution. Performance evaluation is carried out in Section VI. Section VII concludes the paper.

II. RELATED WORK

In the presence of multiple cloudlets, resource management becomes extremely important as it directly impacts edge service quality and system efficiency. Xu et al. [4] formulated a capacitated cloudlet placement problem to minimize the average transmission delay between users and cloudlets and proposed an approximation algorithm to solve it. Jia et al. [5] studied the load balancing problem among multiple cloudlets. Bhatta and Mashayekhy [6] proposed a heuristic cost-aware cloudlet placement approach that guarantees minimum latency for edge services. Wang et al. [7] formulated the dynamic resource allocation problem in edge computing considering user mobility and proposed an online algorithm to solve it by decoupling the problem into a series of solvable sub-problems. However, none of these studies considers the selfish behavior of the users.

Game theory has been widely used to model and analyze different allocation problems. Algorithmic mechanism design deals with efficiently-computable algorithmic solutions in the presence of strategic players who may misreport their input, and it has been used in distributed computing [8]–[10]. Zavadovski et al. [11] proposed an incentive compatible double auction mechanism, called DeCloud, to offer pay-as-you-go edge services, where ad hoc clouds can be spontaneously formed on the edge of the network. Kiani and Ansari [12] proposed a revenue-maximizing auction-based mechanism for edge computing resources. However, the mechanism is not incentive compatible. Ma et al. [13] modeled the resource allocation problem as a three-sided cyclic game (3CG), where edge nodes and service providers cooperate for completing user requests and compete for their own interest. 3CG is proved to have pure-strategy Nash equilibria and an approximation ratio.

Nevertheless, none of the existing work jointly addresses the AP allocation and service placement problem along with determining service pricing in the EC system. In this paper, we propose an online incentive-compatible computation offloading mechanism to address this problem.

III. EDGE COMPUTING SYSTEM MODEL

We consider an EC system with a set of cloudlets, each of which is equipped with an AP (e.g., base station or WiFi hotspots) and edge servers, to provide edge services for users (Fig. 1). A regional cloudlet (or a group of cloudlets) can act as the EC coordinator with the responsibility of collecting system information such as the user requests and system status. We denote a set of cloudlets by $\mathcal{M} = \{1, 2, \dots, M\}$ and a set of users by $\mathcal{N} = \{1, 2, \dots, N\}$. Users join and leave the system dynamically. Each cloudlet $j \in \mathcal{M}$ has one or multiple edge servers with computation capability $\mathbf{F}_j(\tau)$ (i.e., CPU cycles per second) and memory capacity $\mathbf{D}_j(\tau)$ at time τ . Each AP $i \in \mathcal{M}$ can provide service to \mathbf{P}_i users simultaneously and has a bandwidth $\mathbf{B}_i(\tau)$. The cloudlets are interconnected by a wired network (e.g., wide-area network (WAN) or local-area network (LAN)).

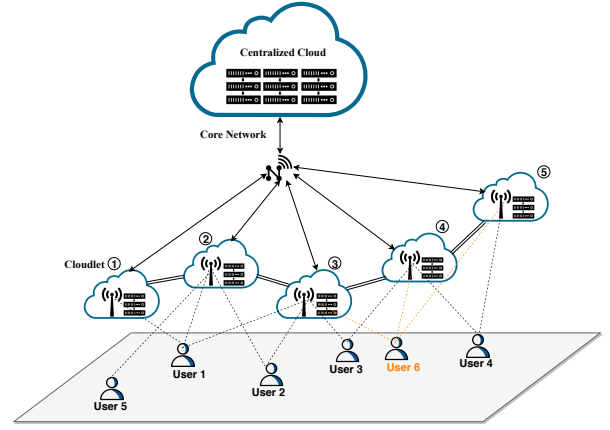


Fig. 1: EC system.

Each user has a computational task requiring remote execution (a user can have multiple tasks, and each is treated independently in this system). The task of user $k \in \mathcal{N}$ is defined by (C_k, D_k, T_k) , where C_k represents the total amount of computational cycles required to obtain the outcome of the task, D_k denotes the data size of the task, and T_k is the maximum tolerable end-to-end latency, measured in time units, for completing the task. Each user can be connected to a cloudlet via an AP through a wireless communication (e.g., WiFi, 4G, or 5G) to offload a task.

A decision pair (i, j) is made by the coordinator for each new joining user k , where $i \in \mathcal{M}$ represents the AP to connect to and $j \in \mathcal{M}$ denotes the assigned edge server at cloudlet $j \in \mathcal{M}$. Even though a user is connected to its nearby AP, its allocated edge server can be at any cloudlet in the EC system. If assigned AP i and edge server j of user k are not associated with each other (i.e., not in the same cloudlet), the system transfers its task from cloudlet i to cloudlet j .

The system state is represented by $\mathcal{I}(\tau) = (P(\tau), Q(\tau))$ at any time instant τ , where $P(\tau)$ and $Q(\tau)$ represent the status of the system in terms of users connected to all APs and computational tasks served by all edge servers at τ , respectively. Specifically, at any time τ , they present the sets of decision variables defined as follows:

$$p_i^k(\tau) = \begin{cases} 1 & \text{if user } k \text{ is connected to AP } i, \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

$$q_j^k(\tau) = \begin{cases} 1 & \text{if user } k \text{ is served by cloudlet } j, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Therefore, at time τ , the total number of users connected to AP i is $u_i(\tau) = \sum_{k \in \mathcal{N}} p_i^k(\tau)$, the total number of computational tasks of users served by cloudlet j is $v_j(\tau) = \sum_{k \in \mathcal{N}} q_j^k(\tau)$, and the total number of computational tasks of users sent to cloudlet j via AP i is $x_{(i,j)}(\tau) = \sum_{k \in \mathcal{N}} p_i^k(\tau) q_j^k(\tau)$. To make the mathematical formulation a linear convex, we can linearize $x_{(i,j)}(\tau)$. We first define a binary decision variable y_{ij} , and define the following set of constraints:

$$y_{ij}(\tau) \geq p_i^k(\tau) + q_j^k(\tau) - 1, \quad \forall i, j \in \mathcal{M}. \quad (3)$$

to ensure that $y_{ij}(\tau)$ is one if both $p_i^k(\tau)$ and $q_j^k(\tau)$ are one; and zero otherwise. We then define:

$$x_{(i,j)}(\tau) = \sum_{i \in \mathcal{M}} \sum_{j \in \mathcal{M}} y_{ij}(\tau). \quad (4)$$

When a new user joins the system, it will impact the system and all existing users. Users using the same AP, edge server, or both as the new user may experience an additional delay. We model the new system state $\tilde{\mathcal{I}}$ after a decision pair (i^*, j^*) is assigned to a new joining user k at time τ , assuming existing users in the system are following their assigned decision pairs. We simply simulate the new system state $\tilde{\mathcal{I}} = (\tilde{P}, \tilde{Q})$ considering $\hat{p}_{i^*}^k(\tau) = 1$ and $\hat{q}_{j^*}^k(\tau) = 1$ while the states of other APs $i \neq i^*$ and edge servers $j \neq j^*$ remain unchanged. In addition, $\hat{u}_{i^*}(\tau) = u_{i^*}(\tau - 1) + 1$ and $\hat{v}_{j^*}(\tau) = v_{j^*}(\tau - 1) + 1$.

When a user leaves the system, the coordinator updates the system state by releasing the communication and computing resources allocated to that user. Specifically, considering an assigned decision pair (i^*, j^*) , we have $\hat{u}_{i^*}(\tau) = u_{i^*}(\tau - 1) - 1$, $\hat{v}_{j^*}(\tau) = v_{j^*}(\tau - 1) - 1$, and other related parameters will be updated, accordingly.

We next describe the AP allocation model, the service placement model, and the end-to-end latency model in detail.

A. Access Point Allocation Model

As mentioned, each cloudlet is associated with an AP. It is possible that a user is within a range of multiple APs and can access any of them, but the user will be connected to only one AP for each of its tasks. We define an indicator variable δ_{ki} that characterizes the availability of AP $i \in \mathcal{M}$ to user $k \in \mathcal{N}$ at time τ as follows:

$$\delta_{ki}(\tau) = \begin{cases} 1 & \text{if AP } i \text{ is available to user } k \text{ at } \tau, \\ 0 & \text{otherwise.} \end{cases}$$

This indicates whether user k can connect to AP i at time τ or not. Therefore, we have the following constraint for AP selection:

$$\sum_{i \in \mathcal{M}} \delta_{ki}(\tau) p_i^k(\tau) = 1, \forall k \in \mathcal{N}, \quad (5)$$

which implies that each user can only connect to one available AP at time τ .

If too many users choose to connect to the same AP simultaneously, they may incur severe interference, which eventually leads to lower uplink data rate. This would negatively affect the performance of computation offloading in the EC system. Therefore, the system needs to guarantee the following:

$$u_i(\tau) \leq \mathbf{P}_i, \forall i \in \mathcal{M}. \quad (6)$$

B. Service Placement Model

When a user requests an edge service, the EC coordinator needs to decide where to properly place computational resources (e.g., VM or Container) to serve this user. Specifically, the requested resources can be hosted on any cloudlet $j \in \mathcal{M}$ that satisfies the QoE requirements of the user and improves the efficiency of the EC system. If tasks can only be executed

on edge servers associated with their connected APs, as the number of arriving users increases, the EC system will be overloaded quickly leading to unsatisfactory performance. Therefore, we seek to find a proper service placement for each user's task.

Each user's task is served by only one cloudlet, thus we have:

$$\sum_{j \in \mathcal{M}} q_j^k(\tau) = 1, \forall k \in \mathcal{N}. \quad (7)$$

In addition, the assignment of tasks to edge servers of each cloudlet should not exceed its capacity:

$$\sum_{k \in \mathcal{N}} q_j^k(\tau) D_k \leq \mathbf{D}_j(\tau), \forall j \in \mathcal{M}. \quad (8)$$

Moreover, we need to ensure that the total number of users connecting to APs is exactly equal to the total number of users served by the cloudlets all the time. Therefore, we have:

$$\sum_{i \in \mathcal{M}} u_i(\tau) = \sum_{j \in \mathcal{M}} v_j(\tau). \quad (9)$$

C. End-to-End Latency Model

End-to-end latency includes the network delay of transmitting the data to a cloudlet (communication delay), the processing time at the cloudlet (computation delay), and finally the network transport delay of transmitting the results to the user's device (communication delay).

C.1) Communication Delay. The communication delay consists of the transmission delay of the user connecting to a proper AP and the transferring delay of the AP relaying to a proper edge server if the connected AP and edge server are not associated with each other.

Transmission Delay. Transmission delay is determined by the wireless communication conditions (e.g., the number of users connected to same AP). Assuming the bandwidth of an AP is equally allocated to all users connecting to it, the bandwidth allocated to user k at time τ from AP i is $r_{ki}(\tau) = \mathbf{B}_i / u_i(\tau)$. Therefore, the uplink transmission delay of offloading task k to AP i at time τ is calculated as:

$$\Lambda_{k,i}^t(\tau) = \frac{D_k}{r_{ki}(\tau)} = \frac{D_k u_i(\tau)}{\mathbf{B}_i}. \quad (10)$$

Transferring Delay. When the connected AP i is not associated with the assigned edge server j , i.e., $i \neq j$, we consider a transferring delay $\Lambda_{i,j}^f(\tau)$ as a function of hop distance between the cloudlet of the connected AP and the desired cloudlet. This is due to the fact that the cloudlets are interconnected via LAN and their physical distance is small. Obviously, if $i = j$, there is no transferring delay, i.e., $\Lambda_{i,j}^f(\tau) = 0$.

Similar to many studies (e.g., [1], [14]), we neglect the delay from the edge server to send the computational results back to the user when the connected AP is associated with the assigned edge server (i.e., no transferring delay). Otherwise, we consider the transferring delay as the total backhaul delay. This is because that the size of computation outcome for many

applications or computational tasks (e.g., image recognition) is usually much smaller than the size of input data.

C.2) Computation Delay. We consider the computational capabilities of a cloudlet are fairly divided among its assigned tasks. The computation delay of the task of user k executed on an edge server of cloudlet j at time τ is calculated using:

$$\Lambda_{k,j}^c(\tau) = \frac{C_k v_j(\tau)}{\mathbf{F}_j}, \forall j \in \mathcal{M}. \quad (11)$$

C.3) Total Delay. The end-to-end latency (total delay) experienced by user k with an assigned decision pair (i, j) at time τ is as follows:

$$\Lambda_{k,(i,j)}^l(\tau) = \Lambda_{k,i}^t(\tau) + 2\Lambda_{i,j}^f + \Lambda_{k,j}^c(\tau) \quad (12)$$

Furthermore, the system needs to ensure that the total delay experienced by user k does not exceed its maximum tolerable end-to-end latency, that is:

$$\Lambda_{k,(i,j)}^l(\tau) \leq T_k, \forall k \in \mathcal{N}. \quad (13)$$

IV. MECHANISM DESIGN-BASED OFFLOADING

Users can be modeled as selfish players that can game the system leading to network congestion, imbalance load, and inefficient resource allocation. Algorithmic mechanism design provides a suitable approach to incentivize players to cooperate with the system in order to reach desirable outcomes. The goal of algorithmic mechanism design is to design a system for such self-interested players, such that their strategies at equilibrium lead to expected system performance. In this section, we propose a computation offloading and pricing mechanism (COPM) to solve the dynamic computation offloading problem in edge computing based on algorithmic mechanism design.

A. Utility Functions

A.1) User-Centric Model. A user $k \in \mathcal{N}$ sends its offloading request in the form of (C_k, D_k, T_k) at time τ to the EC system. The valuation of user k for a decision pair (i, j) considering $\Lambda_{k,(i,j)}^l(\tau)$ (experienced latency) and T_k is defined as:

$$V_{k,T_k}^{(i,j)}(\tau) = \psi_k(T_k - \Lambda_{k,(i,j)}^l(\tau)), \quad (14)$$

where ψ_k is a constant value representing user k 's monetary preference per unit of time for its QoE.

The *utility of user k* when it follows assigned decision pair (i, j) at time τ is determined by:

$$U_{k,T_k}^{(i,j)}(\tau) = \underbrace{V_{k,T_k}^{(i,j)}(\tau)}_{\text{valuation}} - \underbrace{w_k^{(i,j)}(\tau)}_{\text{payment}}, \quad (15)$$

where $w_k^{(i,j)}(\tau)$ is the payment of the user for completing its task through the assigned decision pair. We assume that users are risk-neutral and want to maximize their utilities.

A.2) System-Centric Model. The EC system aims to maximize the social surplus of all current users (excluding new users) while satisfying the QoE of each user.

When new user k joins the system with an assigned decision pair (i^*, j^*) at time τ , we define the *valuation of the system* as follows:

$$V_{s,(i^*,j^*)}(\tau) = \sum_{n \in \mathcal{N} \setminus k} \sum_{i \in \mathcal{M}} \sum_{j \in \mathcal{M}} A \cdot \Lambda_{n,(i,j)}^l(\tau), \quad (16)$$

where A is a vector of parameters α, β, γ representing the monetary preferences per unit of time in transmission, transferring, and computation parts of the offloading, respectively. In particular, α_i is the monetary value of time for AP $i \in \mathcal{M}$; $\beta_{(i,j)}$ is the monetary value of time for transferring a task to edge server j via assigned AP i ; and γ_j is the monetary value of time for edge server $j \in \mathcal{M}$. Therefore, the *valuation of the system* is calculated as:

$$V_{s,(i^*,j^*)}(\tau) = \left(\sum_{i \in \mathcal{M}} \sum_{n \in \mathcal{N} \setminus k} \alpha_i \hat{\Lambda}_{n,i}^t(\tau) + \sum_{i \in \mathcal{M}} \sum_{j \in \mathcal{M}} \beta_{(i,j)} x_{(i,j)}(\tau) 2\hat{\Lambda}_{i,j}^f + \sum_{j \in \mathcal{M}} \sum_{n \in \mathcal{N} \setminus k} \gamma_j \hat{\Lambda}_{n,j}^c(\tau) \right), \quad (17)$$

where $\hat{\Lambda}_{n,i}^t(\tau)$ and $\hat{\Lambda}_{n,j}^c(\tau)$ are the new transmission delay and the new computation delay of the current users ($\forall n \in \mathcal{N} \setminus k$) after user k joins at time τ , respectively. In particular, for all AP $i \neq i^*$ and edge server $j \neq j^*$, we have $\hat{\Lambda}_{n,i}^t(\tau) = \Lambda_{n,i}^t(\tau - 1)$ and $\hat{\Lambda}_{n,j}^c(\tau) = \Lambda_{n,j}^c(\tau - 1)$ for all users. For AP i^* and edge server j^* , the value of $\hat{\Lambda}_{n,i^*}^t(\tau)$ and $\hat{\Lambda}_{n,j^*}^c(\tau)$ will be calculated according to new $\hat{u}_{i^*}(\tau) = u_{i^*}(\tau - 1) + 1$ and new $\hat{v}_{j^*}(\tau) = v_{j^*}(\tau - 1) + 1$, respectively. Note that the transferring delay between any two cloudlets will not be affected by the new joining user (i.e., $\hat{\Lambda}_{i,j}^f = \Lambda_{i,j}^f$), since it depends on their number of hop distances.

The *utility of the system* is defined as:

$$U_{s,(i^*,j^*)}(\tau) = \underbrace{w_s(\tau)}_{\text{payment}} - \underbrace{V_{s,(i^*,j^*)}(\tau)}_{\text{valuation}} \quad (18)$$

Note that the EC system is better off as the total delay of all current users decreases. Moreover, the mechanism is budget balanced, where the exchanged payments are equal. Meaning that: $w_s(\tau) = w_k^{(i,j)}(\tau)$.

B. COPM: Computation Offloading and Pricing Mechanism

The state of the EC system changes dynamically over time by the arrival and departure of users. At any time a new user requests to join, we define a *game* between the new user as a player and all current users in the EC system. The objective is to maximize the utilities of both the new user and the EC system. We propose the *instant social surplus* to handle the dynamic changes of the EC system:

Definition 1 (Instant social surplus). *The instant social surplus at any time τ is the sum of the utility of new user and the utility of the system:*

$$U_{s,(i^*,j^*)}(\tau) + U_{k,T_k}^{(i,j)}(\tau) = V_{k,T_k}^{(i,j)}(\tau) - V_{s,(i^*,j^*)}(\tau). \quad (20)$$

The EC system aims to assign user k a proper decision pair (i, j) maximizing the instant social surplus while guaranteeing this user's QoE and the system's capacity constraints.

Users can choose to lie about their true preferences (i.e., maximum tolerable end-to-end latency) in order to increase their own utility. Such an action could inversely decrease the overall system efficiency. Therefore, designing an *incentive-compatible* mechanism in which users have no incentive to lie about their true preferences is extremely crucial in reality. In an incentive-compatible mechanism, truth-telling is a dominant strategy. As a result, it never pays off for any user to deviate from reporting its true preference, irrespective of what the other users report as their preferences.

We propose an optimal incentive-compatible offloading mechanism, COPM, that consists of a decision pair allocation scheme and a payment determination scheme. To achieve incentive compatibility, we need to design an optimal decision pair allocation scheme (subsection B.1) along with a payment function (subsection B.2) designed based on Vickrey-Clarke-Groves (VCG) pricing [15]. We describe our offloading mechanism design in detail in the following.

B.1) Decision Pair Allocation Scheme. The goal of the EC system is to allocate an optimal decision pair to each new joining user in order to maximize the instant social surplus while satisfying the user's preference. We define the Maximization of Instant Social Surplus problem, called MISS, as follows:

$$\max_{(i,j) \in \mathcal{P}} -V_{s,(i,j)}(\tau) + V_{k,T_k}^{(i,j)}(\tau) \quad (21)$$

s.t. (3)(4)(5)(6)(7)(8)(9)(13) and
integrality for the decision variables,

where \mathcal{P} represents the set of all feasible decision pairs for user k .

Since the MISS optimization problem is complicated, we use the factorization techniques to obtain a simplified version of MISS, called MISS².

Observation 1. *The MISS problem is equivalent to finding a decision pair that minimizes the sum of increase in the total*

delay of all current users (not including the new user) and the total delay of the new user itself.

According to Eq. (14) and (17), the objective function of MISS in Eq. (21) can be rewritten as Eq. (19). The detailed proof is given in our technique report [16].

Since the EC system knows the state of the system, it is easy to find that in Eq. (19), Term 1 representing the utility of the system before new user k joins (i.e., $V_{s,(i,j)}(\tau)$) is constant. Term 2 is also constant as ψ_k and T_k do not depend on the decision pair. Since the values of α_i , γ_j , and ψ_k are predefined constants, the MISS problem has an equivalent minimization problem defined as MISS² as follows:

$$\begin{aligned} \min_{(i,j) \in \mathcal{P}} \left\{ \underbrace{\alpha_i \sum_{n \in \mathcal{N} \setminus k} (\hat{\Lambda}_{n,i}^t(\tau) - \Lambda_{n,i}^t(\tau - 1))}_{\text{increase in transmission delay}} + \right. \\ \left. \underbrace{\gamma_j \sum_{n \in \mathcal{N} \setminus k} (\hat{\Lambda}_{n,j}^c(\tau) - \Lambda_{n,j}^c(\tau - 1))}_{\text{increase in computation delay}} + \right. \\ \left. \underbrace{\psi_k (\hat{\Lambda}_{k,i}^t(\tau) + 2\hat{\Lambda}_{(i,j)}^f + \hat{\Lambda}_{k,j}^c(\tau))}_{\text{total delay of new user}} \right\} \quad (22) \\ \text{s.t. (3)(4)(5)(6)(7)(8)(9)(13) and} \\ \text{integrality for the decision variables.} \end{aligned}$$

Therefore, the objective now becomes to find an optimal decision pair (i^*, j^*) for user k such that the sum of the increase in total delay of all current users (not including user k) after user k joins (first two terms of Eq. (22) in MISS²) and the total delay of user k is minimized (third term of MISS²). After COPM calculates a proper decision pair for each new user by solving MISS², it calculates a corresponding payment for each user.

B.2) Payment Determination Scheme. After solving the MISS² problem, an optimal decision pair (i^*, j^*) is calculated by the EC coordinator for each new joined user k . The coordinator then needs to compute their payments (e.g., each user k should

$$\begin{aligned} & -V_{s,(i,j)}(\tau) + V_{k,T_k}^{(i,j)}(\tau) \\ & = - \left\{ \underbrace{\sum_{i \in \mathcal{M}} \sum_{n \in \mathcal{N} \setminus k} \alpha_i \Lambda_{n,i}^t(\tau - 1) + \sum_{i \in \mathcal{M}} \sum_{j \in \mathcal{M}} 2\beta_{ij} x_{ij}(\tau - 1) \Lambda_{(i,j)}^f + \sum_{j \in \mathcal{M}} \sum_{n \in \mathcal{N} \setminus k} \gamma_j \Lambda_{n,j}^c(\tau - 1)}_{\text{Term 1}} \right\} + \underbrace{\psi_k T_k}_{\text{Term 2}} \\ & - \left\{ \underbrace{\left\{ \alpha_{i^*} \sum_{n \in \mathcal{N} \setminus k} (\hat{\Lambda}_{n,i^*}^t(\tau) - \Lambda_{n,i^*}^t(\tau - 1)) + 2\beta_{i^*j^*} x_{i^*j^*}(\tau - 1) (\hat{\Lambda}_{(i^*,j^*)}^f - \Lambda_{(i^*,j^*)}^f) + \gamma_{j^*} \sum_{n \in \mathcal{N} \setminus k} (\hat{\Lambda}_{n,j^*}^c(\tau) - \Lambda_{n,j^*}^c(\tau - 1)) \right\}}_{\substack{\text{the value of increase in total delay of the current users} \\ \text{(excluding the new user) when new user } k \text{ joins}}} \right\} \\ & + \left\{ \underbrace{\psi_k (\hat{\Lambda}_{k,i^*}^t(\tau) + 2\hat{\Lambda}_{(i^*,j^*)}^f + \hat{\Lambda}_{k,j^*}^c(\tau))}_{\text{the value of total delay of new user } k} \right\} \end{aligned} \quad (19)$$

pay for using its assigned AP i^* and edge server j^*). We define the payment based on the *marginal cost pricing* as follows:

$$w_k^{(i^*, j^*)}(\tau) = V_{s, (i^*, j^*)}(\tau) - V_{s, (i, j)}(\tau - 1) = \left\{ \alpha_{i^*} \sum_{n \in \mathcal{N} \setminus k} (\hat{\Lambda}_{n, i^*}^t(\tau) - \Lambda_{n, i^*}^t(\tau - 1)) + \gamma_{j^*} \sum_{n \in \mathcal{N} \setminus k} (\hat{\Lambda}_{n, j^*}^c(\tau) - \Lambda_{n, j^*}^c(\tau - 1)) \right\}, \quad (23)$$

where $V_{s, (i, j)}(\tau - 1)$ represents the valuation of the EC system right before user k joins, and $V_{s, (i^*, j^*)}(\tau)$ denotes the valuation of the EC system of all current users (not including the joined user k), calculated according to Eq. (17). Considering the first two terms of Eq. (22), the payment of each user k is exactly equal to the increase in valuations of other current users in the EC system.

Our proposed mechanism, COPM, is incentive compatible. Due to the space limitation, the detailed proof of incentive-compatibility of COPM is given in our technique report [16].

V. DAPA: ONLINE ALGORITHMIC-BASED OFFLOADING

We now describe our online algorithmic solution for our mechanism by proposing Dynamic Allocation and Pricing Algorithm (DAPA), presented in Algorithm 1. When any new user k requests edge service with (C_k, D_k, T_k) at time τ , DAPA first finds all available APs and available edge servers for this user (lines 3-4). Also, DAPA has the information of all current users whose tasks are not yet completed when the new user joins (this is updated based on users leaving the system). This information is in $\mathcal{S}(\tau - 1)$, and for each existing user n it consists of allocated AP $n.i$, allocated edge server $n.j$, start time $n.st$, end time $n.et$, complete time, transmission delay, transferring delay, and computation delay. If the arrival time of user k is larger than the end time of any existing user in $\mathcal{S}(\tau - 1)$, it indicates that these users have completed their tasks and left the system before user k joins. Thus, DAPA applies UPDATE() function to update the system state at τ by updating $\hat{u}_i(\tau) \leftarrow u_i(\tau - 1) - 1$, $\hat{v}_j(\tau) \leftarrow v_j(\tau - 1) - 1$, $\hat{\mathbf{D}}_j(\tau) \leftarrow \mathbf{D}_j(\tau - 1) + D_n$, and other related parameters. The information of the completed user n will be removed from $\mathcal{S}(\tau)$ (lines 7-10).

DAPA defines a 2-D array \mathcal{V} and finds the value of instant social surplus by calculating Eq. (22) for each feasible decision pair (lines 11-13). Note that here, DAPA is not solving the MISS or MISS² problem to find the optimal decision pair, but simply calculating the value of Eq. (22) having a decision pair (i, j) . The optimal decision pair (i^*, j^*) with the minimum value is obtained from \mathcal{V} (line 14).

DAPA uses the COMPUTENEWDELAY() function to check if the reported maximum tolerable end-to-end latency of user k can be met. Specifically, the calculated (i^*, j^*) is temporarily assigned to user k and then its total delay $\Lambda_{k, (i, j)}^l(\tau)$ is computed. If $\Lambda_{k, (i^*, j^*)}^l(\tau) \leq T_k$, it implies assigning (i^*, j^*) to user k is feasible (lines 15-18) and the corresponding price for using this pair is calculated using Eq. (23). Otherwise, the

Algorithm 1 DAPA: Dynamic Allocation and Pricing Algorithm for Offloading

- 1: **Input:** User k edge service request: C_k, D_k, T_k
 - 2: **Input:** System state $\mathcal{I}(\tau - 1) = (P(\tau - 1), Q(\tau - 1))$
 - 3: $H^a \leftarrow$ feasible APs for user k
 - 4: $H^e \leftarrow$ feasible edge servers for user k
 - 5: $\mathcal{S}(\tau - 1) \leftarrow$ the set of information of current users
 - 6: $\mathcal{V} \leftarrow \emptyset$ /*2D array of instant social surplus values (22)*/
 - 7: **for** each current user $n \in \mathcal{S}(\tau - 1)$ **do**
 - 8: **if** $k.st > n.et$ **then**
 - 9: $\hat{I}(\tau) \leftarrow$ UPDATE()
 - 10: $\mathcal{S}(\tau) \leftarrow \mathcal{S}(\tau - 1) \setminus n$
 - 11: **for** each AP $i \in H^a$ **do**
 - 12: **for** each edge server $j \in H^e$ **do**
 - 13: $\mathcal{V}[i][j] \leftarrow$ value of Eq. (22)
 - 14: $(i^*, j^*) \leftarrow \arg \min(\mathcal{V})$
 - 15: $\Lambda_{k, (i^*, j^*)}^l(\tau) \leftarrow$ COMPUTENEWDELAY($k, (i^*, j^*)$) based on Eq. (12)
 - 16: **if** $\Lambda_{k, (i^*, j^*)}^l(\tau) \leq T_k$ **then**
 - 17: $w^* \leftarrow$ payment for using (i^*, j^*) based on Eq. (23)
 - 18: **return** $(i^*, j^*), w^*$
 - 19: **else**
 - 20: user k 's request is sent to the cloud
-

request of user k cannot be served by the EC system, and it will be forwarded to the cloud (line 20).

VI. EXPERIMENTAL RESULTS

A. Experimental Setup

1) *EC System Data:* The simulated area is a 500×500 m² square covered by 8 cloudlets, deployed evenly in this area. The effective radius r_i of coverage of each AP i is randomly selected from [75, 100, 125] meters in order to generate the values of the indicator variable. The coverage areas of cloudlets can overlap, which indicates each arriving user may have multiple APs to connect to based on its coordinates. We set the bandwidth of APs obeys Gaussian distribution with mean $\mu = 100$ Mbps and standard deviation $\sigma = 0.25\mu$. The maximum number of users to be served simultaneously by AP i (\mathbf{P}_i) is uniformly selected from [10, 30]. The edge servers are heterogeneous, and each edge server can be equipped with multiple CPU cores. The computation capability of edge servers (\mathbf{F}_j) is uniformly selected from [5, 10] GHz. The transferring delay between two cloudlets is uniformly distributed in [0.1, 0.5] sec. The memory capacity \mathbf{D}_j of each edge server j is 8 GB.

2) *User Data:* The Poisson process plays an important role in modeling systems, as it is usually used in scenarios where the goal is to count the occurrence of certain events happening at a certain rate but completely at random [17]. In this paper, we assume that user arrival events can be modeled as a Poisson process with rate $\lambda = n_a/3600$, where $n_a = 1200$ represents the number of users arriving in the EC system within one hour. Each user k has a computation offloading request, and its location is arbitrary. The data size D_k of user k

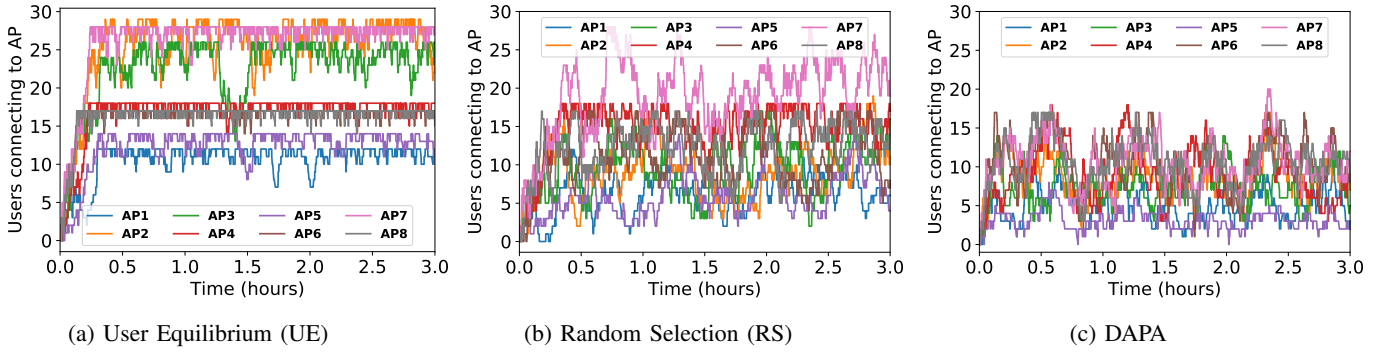


Fig. 2: Analysis of workload on APs

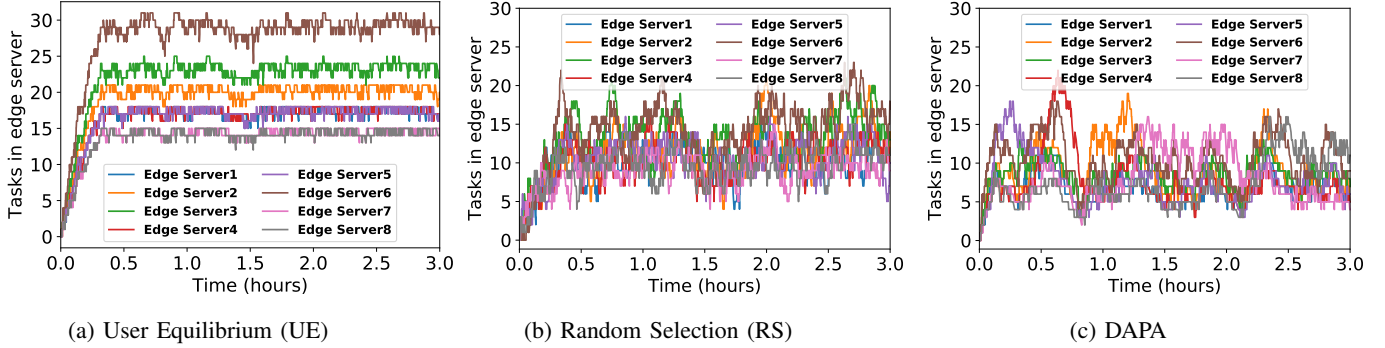


Fig. 3: Analysis of workload on edge servers

is uniformly selected from $[5, 60]$ MB. To specify the required cycles C_k of the computational task, we consider the general application type in which 1 bit requires 1000 cycles to be processed [18]. We roughly classify the users' tasks into three categories: urgent ($t_k + 100$ sec), mid-urgent ($t_k + 200$ sec), and nonurgent ($t_k + 300$ sec), where t_k is the minimum total latency for completing the computational task of user k . We assume that the reported T_k from user k must be no less than t_k . Moreover, ψ_k is 1\$/h, α_i is 50\$/h, and γ_j is 50\$/h.

B. Performance of Benchmark

We simulate a real-time scenario with a duration of 3 hours. To evaluate the performance of our proposed mechanism, DAPA, we compare it with two other offloading strategies:

- 1) **User Equilibrium (UE):** every new user selfishly chooses the decision pair with the minimum total delay.
- 2) **Random Selection (RS):** every new user randomly chooses a feasible decision pair.

We first show the performance of these mechanisms in terms of the workload on APs (Fig. 2) and edge servers (Fig. 3). In particular, these figures show that the dynamics of the number of users $u_i(\tau)$ connecting to each AP i and the number of computational tasks $v_j(\tau)$ on each edge server j over time. The results show that DAPA achieves a more efficient allocation to users such that the workload on each AP and each edge server are balanced overall. Note that RS should be load balanced since it randomly selects a decision pair, however, the experienced time of users by RS is poor (Fig. 4b).

We then investigate the end-to-end latency for completing the computational task of each user when following its assigned decision pair over time. We define T_k^{min} as the minimum latency that the EC system can provide for completing the computational task of user k , that is equal to the value of experienced latency by choosing the UE strategy if only this user exists in the system. We normalize the experienced latency of user k by DAPA (Λ_k^l) and the minimum latency (T_k^{min}) by dividing them by the reported maximum tolerable latency T_k of user k . These normalized values are shown in Fig. 4a. The results show that the experienced end-to-end latency Λ_k^l of user k is different but close to the minimum latency (i.e., T_k^{min} green dots). This figure also shows that the users' preferences are satisfied over time since the experienced latency is always less than or equal to the reported maximum tolerable latency.

Moreover, we study the dynamic changes of the experienced end-to-end latency of users over time in Fig. 4b. The results show that the proposed DAPA outperforms UE and RS in terms of the end-to-end latency that users experience for completing their tasks as the number of joined users increases. This is due to the fact that both UE and RS do not have any policy to consider new users' impacts on other existing users in the system. On the contrary, DAPA aims to find the optimal decision pair for each new user with the objective of jointly minimizing the sum of the increase in total delay of all current users (excluding the new user) after the new user joins and the total delay of the new user. UE leads to the

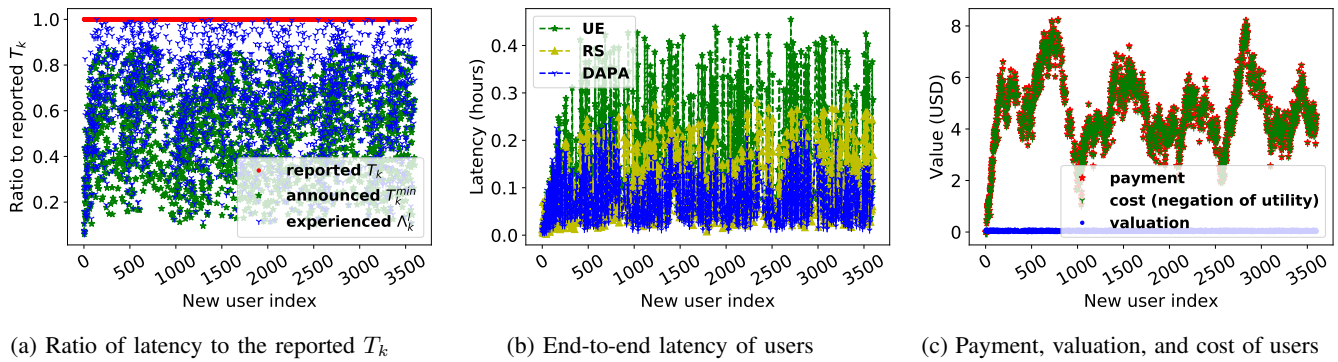


Fig. 4: Analysis of pricing

worst performance as the number of arrived users increases since it considers selfish assignments and the EC system rapidly becomes overloaded on APs/edge servers (as shown in Figs. 2a-3a).

We further evaluate how the EC system makes use of the payments to incentivize each user to report its own true maximum tolerable end-to-end latency. The payments (red points), valuations, and costs (i.e., negative utilities) of joined users are shown in Fig. 4c. The payment of users who join the system at the beginning is much less than users who arrive later. This is due to the fact that each user payment depends on the increase in the end-to-end latency of other existing users in the system (according to Eq. (23)). When there are fewer users, their payment is lower. For example, Figs. 2c and 3c show a decrease in the number of users in the system at 0.6-0.8 hour, that corresponds to about 800th-1000th joining user in Fig. 4c with a reduction in their payments. Also, both of the payments and valuations of the users are always non-negative. Additionally, when all users report their maximum tolerable end-to-end latency truthfully, their costs are minimized (i.e., utilities are maximized) at the equilibrium obtained by DAPA.

VII. CONCLUSION

In this paper, we studied the dynamic computation offloading problem in the EC system. We formulated the computation offloading optimization problem for users joining and leaving the system with the objective of jointly optimizing the access point allocation and service placement problems. To address this challenge, we devised an online incentive-compatible mechanism, DAPA, in which the new users always declare their true preferences. The effectiveness of the mechanism was validated by extensive experiments in comparison to User Equilibrium and Random Selection strategies. For the future work, we plan to consider the effects of user mobility on the computation offloading problem in edge computing.

Acknowledgment. This research was supported in part by NSF grant CNS-1755913.

REFERENCES

- [1] X. Chen, "Decentralized computation offloading game for mobile cloud computing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 26, no. 4, pp. 974–983, 2014.
- [2] M. Satyanarayanan, "The emergence of edge computing," *Computer*, vol. 50, no. 1, pp. 30–39, 2017.
- [3] E. Farhangi Maleki and L. Mashayekhy, "Mobility-aware computation offloading in edge computing using prediction," in *Proc. of the 4th IEEE Intl. Conf. on Fog and Edge Computing*, 2020, pp. 1–6.
- [4] Z. Xu, W. Liang, W. Xu, M. Jia, and S. Guo, "Efficient algorithms for capacitated cloudlet placements," *IEEE Transactions on Parallel and Distributed Systems*, vol. 27, no. 10, pp. 2866–2880, 2015.
- [5] M. Jia, W. Liang, Z. Xu, and M. Huang, "Cloudlet load balancing in wireless metropolitan area networks," in *Proc. of the 35th IEEE Conf. on Computer Communications*, 2016, pp. 1–9.
- [6] D. Bhatta and L. Mashayekhy, "Generalized cost-aware cloudlet placement for vehicular edge computing systems," in *Proc. of the 11th IEEE Intl. Conf. on Cloud Computing Technology and Science*, 2019, pp. 1–8.
- [7] L. Wang, L. Jiao, J. Li, and M. Mühlhäuser, "Online resource allocation for arbitrary user mobility in distributed edge clouds," in *Proc. of the 37th IEEE Intl. Conf. on Dist. Computing Systems*, 2017, pp. 1281–1290.
- [8] L. Mashayekhy, N. Fisher, and D. Grosu, "Truthful mechanisms for competitive reward-based scheduling," *IEEE Transactions on computers*, vol. 65, no. 7, pp. 2299–2312, 2016.
- [9] W. Shi, L. Zhang, C. Wu, Z. Li, and F. Lau, "An online auction framework for dynamic resource provisioning in cloud computing," in *ACM SIGMETRICS Performance Evaluation Review*, vol. 42, no. 1, 2014, pp. 71–83.
- [10] N. Sharghivand, F. Derakhshan, and L. Mashayekhy, "QoS-aware matching of edge computing services to Internet of Things," in *Proc. of the 37th IEEE Intl. Perf. Computing and Comm. Conf.*, 2018, pp. 1–8.
- [11] A. Zavadovski, S. Bayhan, N. Mohan, P. Zhou, W. Wong, and J. Kangasharju, "DeCloud: Truthful decentralized double auction for edge clouds," in *Proc. of the 39th IEEE Intl. Conf. on Distributed Computing Systems*, 2019, pp. 2157–2167.
- [12] A. Kiani and N. Ansari, "Toward hierarchical mobile edge computing: An auction-based profit maximization approach," *IEEE Internet of Things Journal*, vol. 4, no. 6, pp. 2082–2091, 2017.
- [13] S. Ma, S. Guo, K. Wang, W. Jia, and M. Guo, "A cyclic game for joint cooperation and competition of edge resource allocation," in *Proc. of the 39th IEEE Intl. Conf. on Distributed Comp. Systems*, 2019, pp. 503–513.
- [14] X. Ma, C. Lin, X. Xiang, and C. Chen, "Game-theoretic analysis of computation offloading for cloudlet-based mobile cloud computing," in *Proc. of the 18th ACM Intl. Conf. on Modeling, Analysis and Simulation of Wireless and Mobile Systems*, 2015, pp. 271–278.
- [15] N. Nisan, T. Roughgarden, E. Tardos, and V. V. Vazirani, *Algorithmic game theory*. Cambridge university press, 2007.
- [16] W. Ma and L. Mashayekhy, "Truthful Computation Offloading Mechanisms for Edge Computing," *arXiv e-prints*, [Online]. Available: <https://arxiv.org/abs/2006.01553>, 2020.
- [17] D. P. Bertsekas and J. N. Tsitsiklis, *Introduction to probability*. Athena Scientific Belmont, MA, 2002, vol. 1.
- [18] J. Kwak, Y. Kim, J. Lee, and S. Chong, "Dream: Dynamic resource and task allocation for energy minimization in mobile cloud systems," *IEEE Journal on Selected Areas in Communications*, vol. 33, no. 12, pp. 2510–2523, 2015.