

Energy Minimization for Mobile Edge Computing Networks with Time-Sensitive Constraints

Jun-Jie Yu[†], Han Wang[†], Mingxiong Zhao^{*}, Wen-Tao Li, Hui-Qi Bao, Li Yin, and Mi Wu

National Pilot School of Software, Yunnan University, Kunming, China

Email: mx_zhao@ynu.edu.cn

Abstract—Mobile edge computing (MEC) provides users with a high quality experience (QoE) by placing servers with rich services close to the end users. Compared with local computing, MEC can contribute to energy saving, but results in increased communication latency. In this paper, we jointly optimize task offloading and resource allocation to minimize the energy consumption in an orthogonal frequency division multiple access (OFDMA)-based MEC networks, where the time-sensitive tasks can be processed at both local users and MEC server via partial offloading. Since the optimization variables of the problem are strongly coupled, we first decompose the original problem into two subproblems named as offloading selection (P_O), and subcarriers and computing resource allocation (P_S), and then propose an iterative algorithm to deal with them in a sequence. To be specific, we derive the closed-form solution for P_O , and deal with P_S by an alternating way in the dual domain due to its NP-hardness. Simulation results demonstrate that the proposed algorithm outperforms the existing schemes.

Index Terms—MEC, OFDMA, QoE, time-sensitive.

I. INTRODUCTION

With the explosive growth of smart terminals in the IoT era, a large number of applications have emerged that have high computation loads and critical latency (e.g., real-time online games, virtual reality). However, due to the limited energy capacity and computation capability on IoT devices, they cannot fully support the computation-intensive and delay-sensitive services [1]. To find a way out of this dilemma, a novel computing paradigm called Mobile edge computing (MEC) has been envisioned as a promising technology, which integrates cloud computing and mobile network to offer considerable computation resources at network edge to process the offloaded computation intensive or energy-consuming tasks from IoT devices [2].

As an effective method to liberate IoT devices from computation-intensive computing workloads, MEC can efficiently reduce the energy consumption of IoT devices, and thus has been considered as a promising architecture for the scenarios with energy-constrained IoT devices [3], [4]. To be specific, user association was optimized to minimize the total energy consumption of MEC systems in [3]. Meanwhile, the authors in [4] aimed to minimize users' power consumption while trading off the allocated resources for local computation and task offloading.

However, offloading energy-consuming workloads to MEC servers also invokes extra latency which significantly affects

the quality of experience (QoE) of users, and thus cannot be ignored in the system design. As one of the important metrics to measure the performance of MEC network, time sensitivity has triggered more and more research interests and has been investigated in the literature [5]–[7]. To be specific, the authors investigated the latency-minimization problem in a multi-user time-division multiple access MEC offloading system [5]. Moreover, some pioneering works considered the energy consumption minimization with the QoE requirement of time-sensitive computation tasks [6], [7]. To satisfy user demands of various IoT applications, the authors in [6] found a tradeoff between the energy consumption and latency, and formalized the problem into a constrained multi-objective optimization problem. At the same time, the authors in [7] minimized the weighted sum of the execution delay and energy consumption while guaranteeing the transmission power constraint of IoT devices based on partial offloading.

To further improve the utilization of radio resources, Orthogonal Frequency-Division Multiple Access (OFDMA) has been widely applied to the MEC system for various objectives in MEC system, such as profit maximization [8], delay minimization [9], the energy consumption minimization with or without the requirement of computation latency [10], [11], and the energy-latency tradeoff [12]. Specifically, the authors in [8] proposed a priority-based task scheduling policy and jointly optimized the computation and communication resource to maximize the profit of mobile network operator while satisfying users' quality of service (QoS). Furthermore, to minimize the maximum delay of each mobile device, the authors in [9] considered a partial offloading scheme and developed a heuristic algorithm to jointly optimize the subcarrier and power allocation. At the meantime, QoS was introduced to minimize the energy consumption of a multi-cell MEC network in [10]. Moreover, the authors proposed an energy-efficient joint offloading and wireless resource allocation strategy for delay-critical applications to minimize the total (or weighted-sum) energy consumption of the mobile devices [11], and the tradeoff between energy consumption and sensitive latency was further considered to design the energy-aware offloading scheme in [12].

However, the inspiring works [11], [12] only investigated binary offloading without the consideration of partial offloading, which can flexibly allocate resources for computation offloading and local computing, and thus achieve better performance [13], although the QoE requirement of time-sensitive computation tasks was taken into account [12]. Motivated by the aforementioned issues, we consider partial offloading where mobile

[†]Indicates equal contribution. *M. Zhao is the corresponding author. This work is supported in part by the National Natural Science Foundation of China under Grant 61801418, in part by Yunnan Applied Basic Research Projects 2019FD-129.

date can be computed at both local users and MEC server, and investigate an OFDMA-based MEC network in this paper. Our designed framework aims to minimize the total energy consumption of the whole network, where offloading ratio, computation capability and subcarriers are jointly optimized to satisfy the QoE requirement of time-sensitive computation tasks of users.

II. SYSTEM MODEL AND PROBLEM FORMULATION

We consider an OFDMA-based MEC system with K users and one base-station (BS) integrated with an MEC server to execute the offloaded data of users, and all nodes are equipped with a single antenna. Denote $\mathcal{K} \triangleq \{1, 2, \dots, K\}$ as the set of users, and let $\mathcal{N} \triangleq \{1, 2, \dots, N\}$ be the index for multiple orthogonal subcarriers, each of which has bandwidth B and can be assigned to only one user. In this system, we assume that user k has a task described by a tuple of four parameters $\{R_k, c_k, \lambda_k, t_k\}$, where R_k indicates the amount of input data to be processed, c_k represents the number of CPU cycles for computing 1-bit of input data, $\lambda_k \in [0, 1]$ is the proportion of R_k offloading to MEC, while the rest $(1 - \lambda_k)R_k$ bits are processed by its local CPU, and t_k is the maximum tolerable latency. In this paper, it is assumed that the maximum tolerable latency for user k , t_k is no longer than the channel coherence time, such that the wireless channels remain constant during a time slot with length T , i.e., $t_k \leq T, \forall k$, but can vary from time to time. The local CPU frequency of user k is characterized by f_k , and $f_{k,m}$ is the computational speed of the edge cloud assigned to user k , where both of them are measured by the number of CPU cycles per second. Herein, a practical constraint that the total computing resources allocated to all the associated users must not exceed the servers computing capacity F , is given by $\sum_{k \in \mathcal{K}} f_{k,m} \leq F$.

In the following, the time latency and the energy consumption of user k for our considered system are given in details.

A. Latency

1) *Local Computing at Users:* Consider the local computing for executing the residual $(1 - \lambda_k)R_k$ input bits at user k , the time consumption for local computing at user k is

$$t_{k,l} = \frac{c_k (1 - \lambda_k) R_k}{f_k}. \quad (1)$$

2) *Computation Offloading:* According to the OFDMA mechanism, the inter-interference is ignored in virtue of the exclusive subcarrier allocation. Therefore, the aggregated transmission rate to offload $\lambda_k R_k$ input bits from user k to MEC server is expressed as

$$r_k = B \sum_{n \in \mathcal{N}} x_{k,n} \log_2 \left(1 + \frac{p_{k,n} g_{k,n}}{\sigma_n^2} \right), \quad (2)$$

where $g_{k,n}$ and σ_n^2 are the channel gain between user k and BS, and the variance of the additive white Gaussian noise at BS on subcarrier n , respectively, where we set $\sigma_n^2 = \sigma^2, \forall n$. Denote $p_{k,n}$ as the transmission power of user k on subcarrier n . Apparently, any power optimization solution

has a good influence on the system performance. For the sake of simplicity, we set $p_{k,n} = p_k^{\max} / N_k$, where p_k^{\max} is the maximum transmission power and N_k denotes the number of subcarriers allocated to user k . Moreover, if user k does not offload data to the MEC server, $p_{k,n} = 0$. Meanwhile, denote $x_{k,n}$ as the channel allocation indicator, specifically $x_{k,n} = 1$ means that subcarrier n is assigned to user k , otherwise $x_{k,n} = 0$.

The offloading time $t_{k,\text{off}}$ of user k mainly consists of two parts: the uplink transmission time $t_{k,u}$ from user k to MEC server and the corresponding execution time at MEC server $t_{k,m}$. Therefore, the offloading time $t_{k,\text{off}}$ is given by

$$t_{k,\text{off}} = t_{k,u} + t_{k,m} = \frac{\lambda_k R_k}{r_k} + \frac{\lambda_k R_k c_k}{f_{k,m}}. \quad (3)$$

Due to the parallel computing at users and MEC server, the total latency for user k depends on the larger one between $t_{k,l}$ and $t_{k,\text{off}}$, and can be expressed as $t_k = \max\{t_{k,l}, t_{k,\text{off}}\}$.

B. Energy Consumption

According to the strategy of computation offloading, the total energy consumption comprises two parts: the energy for local computing and for offloading, given in details as follow.

1) *Local Computing mode:* Given the processor's computing speed f_k , the power consumption of the processor is modeled as $\kappa_k f_k^3$ (joule per second), where κ_k represents the computation energy efficiency coefficient related to the processor's chip of user k [14]. Taking consideration of (1), the energy consumption at this mode is given by

$$E_{k,l} = \kappa_k f_k^3 t_{k,l} = \kappa_k c_k (1 - \lambda_k) R_k f_k^2. \quad (4)$$

2) *Computation offloading mode:* In this mode, the energy consumption includes the cost of uplink transmitting ($E_{k,u}$) and remote computing for offloaded $\lambda_k R_k$ input bits ($E_{k,m}$), which can be obtained as

$$E_{k,\text{off}} = \sum_{n \in \mathcal{N}} x_{k,n} p_{k,n} \frac{\lambda_k R_k}{r_k} + \kappa_m \lambda_k c_k R_k f_{k,m}^2, \quad (5)$$

where κ_m is the computation energy efficiency coefficient related to the processor's chip of MEC server.

Therefore, the total energy consumption for user k related with its computation offloading strategy in our system is $E_k = E_{k,l} + E_{k,u} + E_{k,m}$.

In this paper, we minimize the overall energy consumption of the considered system as follows

$$\mathbf{P} : \min_{\lambda, \mathbf{f}, \mathbf{X}} \sum_{k \in \mathcal{K}} E_k \quad (6a)$$

$$\text{s.t. } 0 \leq \lambda_k \leq 1, \forall k, \quad (6b)$$

$$\max\{t_{k,l}, t_{k,\text{off}}\} \leq T, \forall k, \quad (6c)$$

$$0 \leq f_{k,m}, \forall k, \quad (6d)$$

$$\sum_{k \in \mathcal{K}} f_{k,m} \leq F, \quad (6e)$$

$$\sum_{k \in \mathcal{K}} x_{k,n} \leq 1, \forall n, \quad (6f)$$

$$x_{k,n} \in \{0, 1\}, \forall k, n, \quad (6g)$$

which is related to resource allocation on subcarriers, offloading communication and computation, and $\lambda \triangleq \{\lambda_k\}$, $\mathbf{f} \triangleq \{f_{k,m}\}$ and $\mathbf{X} \triangleq \{x_{k,n}\}$. The constraints above can be explained as follows: constraint (6c) states that the task of user k must be completely executed within a time slot; constraint (6d) and (6e) show that MEC server must allocate a positive computing resource to the user associated with it, and the sum of which cannot exceed the total computational capability of MEC server; constraint (6f) and (6g) enforce that each subcarrier can only be used by one user to avoid the multi-user interference.

III. PROPOSED ALGORITHM

In this section, we provide offloading and resource allocation strategy for the considered optimization problem \mathbf{P} , which is intractable to deal with due to the coupled variants in both the constraints and the objective function based on our observation. To decouple these variants, we will divide the original problem \mathbf{P} into two subproblems: 1) \mathbf{P}_O , offloading ratio selection; 2) \mathbf{P}_S , subcarriers and computing resource allocation.

Firstly, with given computation capability assignment \mathbf{f} and subcarrier allocation strategy \mathbf{X} , we can obtain the optimal offloading ratio λ^* at the outer loop. **Secondly**, with the newly obtained offloading ratio λ^* , we can optimize (\mathbf{f}, \mathbf{X}) at one iteration at its inner loop, and renew auxiliary variable ϕ . Then, with the newly optimized (\mathbf{f}, \mathbf{X}) , we further update the dual variables (α, β, γ) at the next iteration at its inner loop. **Finally**, we will iteratively update the derived $(\lambda, \mathbf{f}, \mathbf{X})$ at the outer loop, and the procedures are known as the BCD method [15], [16]. In this section, the joint optimization on offloading ratio, and subcarriers and computing resource allocation will be proposed in accordance with the iterative approach based on BCD method as follows.

A. Offloading ration selection

Given computation capability assignment \mathbf{f} and subcarrier allocation strategy \mathbf{X} , the optimal offloading ratio λ^* can be obtained by solving the following problem,

$$\mathbf{P}_O : \min_{\lambda} \sum_{k \in \mathcal{K}} E_k \quad (7)$$

s.t. (6b)(6c),

which can be decoupled into K subproblems for each user, given by

$$\mathbf{P}_{O1} : \min_{\lambda_k} E_k \quad (8a)$$

$$\text{s.t. } 0 \leq \lambda_k \leq 1, \quad (8b)$$

$$1 - \frac{Tf_k}{c_k R_k} \leq \lambda_k \leq \frac{Tr_k f_{k,m}}{R_k f_{k,m} + r_k R_k c_k}, \quad (8c)$$

where (8c) can be derived with the help of (1) and (3). Apparently, \mathbf{P}_{O1} is a convex problem with respect to (w.r.t.) λ_k , and the optimal offloading ratio λ_k^* at user k can be achieved according to the following theorem.

Theorem 1: Given (\mathbf{f}, \mathbf{X}) , the optimal λ_k^* for \mathbf{P}_{O1} is

$$\lambda_k^* = \begin{cases} \max \left\{ 1 - \frac{Tf_k}{c_k R_k}, 0 \right\}, & \text{if } \frac{\partial E_k}{\partial \lambda_k} \geq 0, \\ \min \left\{ \frac{Tr_k f_{k,m}}{R_k f_{k,m} + r_k R_k c_k}, 1 \right\}, & \text{otherwise.} \end{cases} \quad (9)$$

Proof: It can be obtained resorting to the first-order condition and comparing with the boundaries points provided by (8b) and (8c). \square

B. Subcarriers and computing resource allocation strategy

With the newly obtained offloading ratio λ^* , the resource allocation strategy will be designed to assign the computation capability of MEC server, and allocate the subcarriers for each user, to further reduce the energy consumption. In this subsection, we aim to minimize the energy consumption for computation offloading mode via the following optimization problem, the objective function of which is $\sum_{k \in \mathcal{K}} E_{k,\text{off}}$ indeed,

$$\mathbf{P}_S : \min_{\mathbf{f}, \mathbf{X}} \sum_{k \in \mathcal{K}} \sum_{n \in \mathcal{N}} x_{k,n} p_{k,n} \frac{\lambda_k R_k}{r_k} + \kappa_m \lambda_k c_k R_k f_{k,m}^2 \quad (10a)$$

$$\text{s.t. (6d) - (6g).}$$

$$t_{k,u} + t_{k,m} \leq T, \forall k, \quad (10b)$$

where the time-sensitive constraint (6c) can be recast as (10b) when λ is given.

However, we cannot transform the primal domain of \mathbf{P}_S into the dual domain directly since r_k is in the denominator and its form in (2) makes the problem more intractable. Therefore, a new non-negative auxiliary variable $\phi \triangleq \{\phi_k\}$ will be introduced to transform \mathbf{P}_S into the following problem,

$$\mathbf{P}_{S1} : \min_{\mathbf{f}, \mathbf{X}, \phi} \sum_{k \in \mathcal{K}} \sum_{n \in \mathcal{N}} x_{k,n} p_{k,n} \frac{\lambda_k R_k}{\phi_k} + \kappa_m \lambda_k c_k R_k f_{k,m}^2 \quad (11a)$$

$$\text{s.t. (6b) - (6g), (10b),}$$

$$0 \leq \phi_k \leq r_k, \forall k, \quad (11b)$$

where we can update ϕ to help minimize the energy consumption. The Lagrangian for the above problem with the given λ^* can be written as

$$\begin{aligned} \mathcal{L}(\mathbf{f}, \mathbf{X}, \phi, \alpha, \beta, \gamma) = & \sum_{k \in \mathcal{K}} \sum_{n \in \mathcal{N}} x_{k,n} p_{k,n} \frac{\lambda_k R_k}{\phi_k} \\ & + \sum_{k \in \mathcal{K}} \kappa_m \lambda_k c_k R_k f_{k,m}^2 + \sum_{k \in \mathcal{K}} \alpha_k \left(\frac{\lambda_k R_k}{\phi_k} + \frac{\lambda_k R_k c_k}{f_{k,m}} - T \right) \\ & + \sum_{k \in \mathcal{K}} \beta_k (\phi_k - r_k) + \gamma \left(\sum_{k \in \mathcal{K}} f_{k,m} - F \right), \end{aligned} \quad (12)$$

where $\alpha \triangleq \{\alpha_k\}$, $\beta \triangleq \{\beta_k\}$, and γ are the non-negative Lagrange multipliers corresponding to the related constraints. Define \mathcal{F} as all sets of possible \mathbf{f} that satisfy constraint (6d),

\mathcal{X} as all sets of possible \mathbf{X} that satisfy constraints (6f) and (6g). The Lagrange dual function is then defined as

$$\mathbf{g}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \gamma) = \min_{\mathbf{f} \in \mathcal{F}, \mathbf{X} \in \mathcal{X}, \phi \in \mathcal{Q}} \mathcal{L}(\mathbf{f}, \mathbf{X}, \phi, \boldsymbol{\alpha}, \boldsymbol{\beta}, \gamma). \quad (13)$$

Furthermore, the Lagrange dual problem is given by

$$\begin{aligned} \max \quad & \mathbf{g}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \gamma) \\ \text{s.t.} \quad & \boldsymbol{\alpha} \succeq \mathbf{0}, \boldsymbol{\beta} \succeq \mathbf{0}, \gamma \geq 0. \end{aligned} \quad (14)$$

With the given offloading ratio and transmission power, the following steps for updating are adopted to obtain the optimal solutions for computing resource and subcarrier allocation.

1) *Computational capabilities assignment*: Employing the KKT conditions, the following condition is both necessary and sufficient for computation capability assignment's optimality:

$$\frac{\partial \mathcal{L}}{\partial f_{k,m}} = 2f_{k,m} \kappa_m \lambda_k R_k c_k - \frac{\alpha_k c_k R_k \lambda_k}{f_{k,m}^2} + \gamma = 0. \quad (15)$$

However, it is difficult to find a closed-form expression for the optimal solution, $f_{k,m}^*$. Fortunately, we can resort to the following proposition to obtain $f_{k,m}^*$.

Since \mathcal{L} is a convex function of $f_{k,m}$, and $\frac{\partial \mathcal{L}}{\partial f_{k,m}}$ increases monotonically along with $f_{k,m}$, we can adopt the bisection method to obtain $f_{k,m}^*$ within $0 \leq f_{k,m} \leq F$. The detailed process of achieving $f_{k,m}^*$ is summarized in Algorithm 1.

Algorithm 1 Proposed Binary Search Algorithm

Input: Given offloading ratio $\boldsymbol{\lambda}^*$, accuracy indicator ϵ and transmission power p .

- 1: **for** $k \in \mathcal{K}$ **do**
 - 2: **Initialize:** $f_{k,m}^{\text{UB}} = F$ **and** $f_{k,m}^{\text{LB}} = 0$;
 - 3: **repeat**
 - 4: Set $f_{k,m} = \frac{1}{2}(f_{k,m}^{\text{UB}} + f_{k,m}^{\text{LB}})$;
 - 5: Compute $\frac{\partial \mathcal{L}}{\partial f_{k,m}}$ according to (15);
 - 6: **if** $\frac{\partial \mathcal{L}}{\partial f_{k,m}} > 0$ **then**
 - 7: set $f_{k,m}^{\text{UB}} = f_{k,m}$;
 - 8: **else**
 - 9: set $f_{k,m}^{\text{LB}} = f_{k,m}$;
 - 10: **until** $\left| \frac{\partial \mathcal{L}}{\partial f_{k,m}} \right| \leq \epsilon_3$;
 - 11: Obtain the optimal $f_{k,m}^*$.
-

2) *Subcarrier allocation strategy*: With the achieved optimal computing capability assignment, the optimal subcarrier allocation can be obtained through the following procedures. With some mathematic manipulations, we can rewrite (12) as

$$\begin{aligned} \mathcal{L}(\mathbf{X}, \phi, \boldsymbol{\alpha}, \boldsymbol{\beta}, \gamma) = & \sum_{k \in \mathcal{K}} \sum_{n \in \mathcal{N}} x_{k,n} \left[p_{k,n} \frac{\lambda_k R_k}{\phi_k} \right. \\ & \left. - B\beta_k \log_2(1 + p_{k,n} \tilde{g}_{k,n}) \right] + \sum_{k \in \mathcal{K}} \omega_k - \gamma F, \end{aligned} \quad (16)$$

where

$$\begin{aligned} \omega_k = & \alpha_k \left(\frac{\lambda_k R_k}{\phi_k} + \frac{\lambda_k R_k c_k}{f_{k,m}} - T \right) + \gamma f_{k,m} + \beta_k \phi_k \\ & + \kappa_m \lambda_k c_k R_k f_{k,m}^2. \end{aligned} \quad (17)$$

On the observation of (16), we further suppose that subcarrier n is assigned to user k , we have

$$\mathcal{L} = \sum_{n \in \mathcal{N}} \mathcal{L}_n + \sum_{k \in \mathcal{K}} \omega_k - \gamma F, \quad (18)$$

where

$$\mathcal{L}_n = p_{k,n} \frac{\lambda_k R_k}{\phi_k} - B\beta_k \log_2(1 + p_{k,n} \tilde{g}_{k,n}). \quad (19)$$

Thus, the subproblem is given by

$$\min_{\mathbf{X}_n \in \mathcal{X}} \mathcal{L}_n(\mathbf{X}_n, \phi, \boldsymbol{\beta}), \quad (20)$$

where $\mathbf{x}_n = \{x_{k,n}\}_{k=1}^K$, and this problem can be solved independently. To minimize each \mathcal{L}_n , the optimal \mathbf{x}_n can be obtained as

$$x_{k,n}^* = \begin{cases} 1, & \text{if } k = k^* = \operatorname{argmin}_k \mathcal{L}_n, \\ 0, & \text{otherwise.} \end{cases} \quad (21)$$

3) *Auxiliary variable ϕ^* selection*: In this part, with the newly obtained $\{\mathbf{f}^*, \mathbf{X}^*\}$ in the above subsections, we now try to find the optimal ϕ^* , which can be solved via the following optimization problem,

$$\mathbf{P}_{\mathbf{S}2} : \min_{\phi} \sum_{k \in \mathcal{K}} \left[\sum_{n \in \mathcal{N}_k} \frac{p_{k,n} \lambda_k R_k}{\phi_k} + \frac{\alpha_k \lambda_k R_k}{\phi_k} + \beta_k \phi_k \right] \quad (22a)$$

s.t. (10b),

$$0 \leq \phi_k \leq \tilde{r}_k, \forall k, \quad (22b)$$

where $\tilde{r}_k \triangleq \sum_{n \in \mathcal{N}_k} B \log_2(1 + p_{k,n} \tilde{g}_{k,n})$, and \mathcal{N}_k denotes the set of subcarriers allocated to user k through (21). It is obvious that $\mathbf{P}_{\mathbf{S}2}$ can be further decoupled into K subproblems w.r.t. each user, given by

$$\mathbf{P}_{\mathbf{S}3} : \min_{\phi_k} \left[\frac{\lambda_k R_k}{\phi_k} \sum_{n \in \mathcal{N}_k} p_{k,n} + \frac{\alpha_k \lambda_k R_k}{\phi_k} + \beta_k \phi_k \right] \quad (23a)$$

s.t. (22b)

$$\frac{\lambda_k R_k}{\phi_k} + \frac{\lambda_k R_k c_k}{f_{k,m}} \leq T, \quad (23b)$$

the optimal auxiliary variable ϕ^* can be obtained by the following theorem.

Theorem 2: Given the optimal computation capability assignment \mathbf{f}^* , and the optimal subcarrier allocation strategy \mathbf{X}^* , the optimal auxiliary variable $\phi^*, \forall k$ is given by

$$\phi_k^* = \begin{cases} \phi_{k,1}, & \text{if } \phi_k^o < \phi_{k,1}, \\ \phi_k^o, & \text{if } \phi_{k,1} \leq \phi_k^o \leq \tilde{r}_k, \\ \tilde{r}_k, & \text{otherwise,} \end{cases} \quad (24)$$

where $\phi_{k,1}$ and ϕ_k^o are given, respectively, by

$$\phi_{k,1} = \frac{\lambda_k R_k f_{k,m}}{T f_{k,m} - \lambda_k R_k c_k}, \quad (25)$$

$$\phi_k^o = \sqrt{\left(\alpha_k + \sum_{n \in \mathcal{N}_k} p_{k,n} \right) \tilde{\lambda}_k}, \quad (26)$$

where $\tilde{\lambda}_k = \frac{\lambda_k R_k}{\beta_k}$.

Proof: Due to the space limitation, the proof is omitted and presented in [17]. \square

TABLE I: SIMULATION PARAMETERS

MEC System Parameters	Values
The CPU frequency of MEC sever	10GHz
The CPU frequency of mobile users	0.6-0.7GHz
The transmission power of users	27.8dBm
Input data size of users	1000 – 1500 bits
Maximum accomplished deadline T	2ms
The computation workload/intensity c_k	1000 – 1200 cycles/bit
Background noise σ^2	10^{-13} W
Subcarrier bandwidth B	12.5KHz
Lagrange Iteration Parameters	Values
Maximum number of iterations z_{\max}	600
Iterations precision	10^{-5}
Stepsizes ζ_k, η_k, θ	$10^{-18}, 10^{-6}, 10^{-5}$

4) *Lagrange Multipliers Update*: In this subsection, since \mathbf{f}^* , \mathbf{X}^* and ϕ^* are obtained, we can deal with the dual problem (14), which is a convex function, by updating α , β and γ using the subgradient method. The dual variables $(\alpha, \beta, \beta, \gamma)$ can be updated according to the following formulations,

$$\begin{aligned} \alpha_k^{z+1} &= \left[\alpha_k^z - \zeta_k \left(\frac{\lambda_k R_k}{\phi_k} + \frac{\lambda_k R_k c_k}{f_{k,m}} - T \right) \right]^+, \\ \beta_k^{z+1} &= [\beta_k^z - \xi_k (\phi_k - r_k)]^+, \\ \gamma^{z+1} &= \left[\gamma^z - \theta \left(\sum_{k \in \mathcal{K}} f_{k,m} - F \right) \right]^+, \end{aligned} \quad (27)$$

where ζ_k , ξ_k , and θ are the stepsizes related to each dual variable during iterations, given in TABLE I.

According to III-B and III-A, the details to solve \mathbf{P} are summarized in Algorithm 2 as follows.

Algorithm 2 Proposed Algorithm

- 1: **Initialization**: Given $\{\mathbf{f}, \mathbf{X}, \mathbf{P}, \phi, \alpha, \beta, \gamma\}$, we set $z = 0$, and denote z_{\max} as the maximum number of iterations.
- 2: **repeat**(from 2 to 14)
- 3: Determine offloading ratio λ according to (9);
- 4: **repeat**(from 4 to 12)
- 5: **repeat**(from 5 to 10)
- 6: Assign computing resource \mathbf{f} by Algorithm 1;
- 7: Determine subcarrier allocation \mathbf{X} via (21);
- 8: Update $p_{k,n} = \frac{p_k^{max}}{N_k}$;
- 9: Update ϕ^* ;
- 10: **until** Lagrangian function converges;
- 11: Update α , β , and γ resorting to (27);
- 12: **until** α, β, γ converge;
- 13: $z = z + 1$;
- 14: **until** $z > z_{\max}$.

IV. NUMERICAL RESULTS

In the simulations, the path loss model is Rayleigh distributed and denoted by $|\beta|d^{-2}$, where $|\beta|$ and d represent the short-term channel fading and the distance between two nodes, respectively. User devices have the same maximum transmit power, which are evenly and independently distributed in a circular area around the MEC server with a radius of 30 meters. Moreover, we set κ_k and κ_m as 10^{-24} and 10^{-26} , respectively, and other

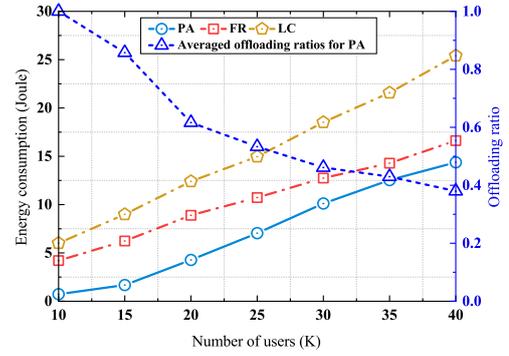


Fig. 1: The comparison of the total energy consumption (and the averaged offloading ratios for PA where $K = 20$) versus different number of users where $N = 512$.

parameters employed in the simulations are summarized in TABLE 1, unless otherwise mentioned.

This section presents the numerical results to demonstrate the better performance of our proposed algorithm compared with the conventional schemes: 1) **Local computation algorithm (LC)**, where all user devices process their own tasks by local CPU. 2) **Fixed offloading ratio algorithm (FR)** where we assign fixed offloading ratios for each user satisfying the time-sensitive constraints (6c).

In Fig.1, for the proposed and reference schemes, we plot the total energy consumption w.r.t. different number of users. With the growth of the number of users, the total energy consumption of different algorithms is increasing, and our proposed algorithm (PA) can help save 40% – 70% and 20%-50% energy consumption compared with LC and FR. Moreover, it can be seen that the averaged offloading ratio will decrease with the increase of the number of users, since computing resource allocated to each user by MEC is declining, and thus users will cut down the offloading ratio.

In Fig.2, we study the influence of the QoE requirements of time-sensitive computation tasks of users on the total energy consumption for different numbers of users. On the observation of Fig.2, the total energy consumption is reducing along with the gradually decreasing QoE requirement of time-sensitive computation tasks. This is because users can ask for more help from MEC server by uploading more data since the requirement of latency is not so strict shown by the blue line, and thus users can reduce more energy dissipation.

In Fig.3, we evaluate the effect of the number of subcarriers on the total energy consumption. Moreover, it can be seen that the total energy consumption is reducing along with the increasing number of subcarriers. The reason is that the user will have a better chance to select the subcarrier with preferable channel gain to get the reduction of transmission power in turn, meanwhile, the user can offload more data with the same amount of transmission power shown by the blue line, and thus reduce the energy consumption with the aid of MEC server.

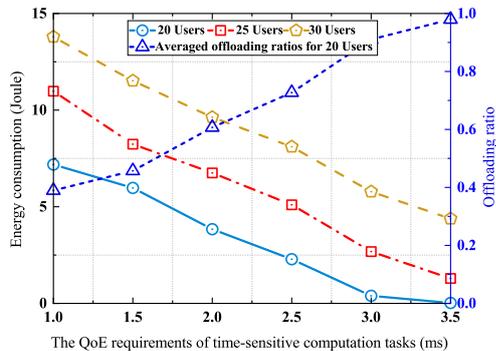


Fig. 2: The comparisons of the total energy consumption (and the averaged offloading ratios where $K=20$) versus different QoE requirements of time-sensitive computation tasks of users where $N=512$.

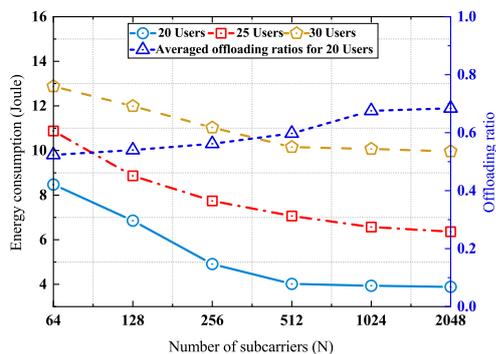


Fig. 3: The comparisons of the total energy consumption versus different number of subcarriers where $K=20, 25, 30$.

V. CONCLUSIONS

In this paper, the energy minimization was investigated by jointly optimizing partial offloading strategy and resource allocation for the OFDMA-based MEC networks, while satisfying the QoE requirements of time-sensitive computation tasks for all users. Due to the coupled optimization variables, we decomposed the formulated minimization problem into two subproblems, and optimized them sequentially instead of solving the original problem directly. Simulation results demonstrate that the proposed algorithm can effectively reduce the total energy consumption of the network, and outperform the reference schemes with a better performance.

REFERENCES

- [1] N. Abbas, Y. Zhang, A. Taherkordi, and T. Skeie, "Mobile edge computing: A survey," *IEEE Internet Things J.*, vol. 5, no. 1, pp. 450–465, 2018.
- [2] S. Mu, Z. Zhong, D. Zhao, and M. Ni, "Joint job partitioning and collaborative computation offloading for Internet of Things," *IEEE Internet Things J.*, vol. 6, no. 1, pp. 1046–1059, 2019.
- [3] Y. Dai, D. Xu, S. Maharjan, and Y. Zhang, "Joint computation offloading and user association in multi-task mobile edge computing," *IEEE Trans. Veh. Tech.*, vol. 67, no. 12, pp. 12 313–12 325, 2018.
- [4] C.-F. Liu, M. Bennis, M. Debbah, and H. V. Poor, "Dynamic task offloading and resource allocation for ultra-reliable low-latency edge computing," *IEEE Trans. Commun.*, vol. 67, no. 6, pp. 4132–4150, 2019.

- [5] J. Ren, G. Yu, Y. Cai, and Y. He, "Latency optimization for resource allocation in mobile-edge computation offloading," *IEEE Trans. Wireless Commun.*, vol. 17, no. 8, pp. 5506–5519, 2018.
- [6] L. Cui, C. Xu, S. Yang, J. Z. Huang, J. Li, X. Wang, Z. Ming, and N. Lu, "Joint optimization of energy consumption and latency in mobile edge computing for Internet of Things," *IEEE Internet Things J.*, vol. 6, no. 3, pp. 4791–4803, 2018.
- [7] Z. Kuang, L. Li, J. Gao, L. Zhao, and A. Liu, "Partial offloading scheduling and power allocation for mobile edge computing systems," *IEEE Internet Things J.*, vol. 6, no. 4, pp. 6774–6785, 2019.
- [8] P. Paymard, N. Mokari, and M. Orooji, "Task Scheduling Based on Priority and Resource Allocation in Multi-User Multi-Task Mobile Edge Computing System," in *Proc. IEEE PIMRC*, 2019, pp. 1–7.
- [9] M. Li and et al, "Joint subcarrier and power allocation for OFDMA based mobile edge computing system."
- [10] A. Khalili and et al, "Joint resource allocation and offloading decision in mobile edge computing," *IEEE Commun. Lett.*, vol. 23, no. 4, pp. 684–687, 2019.
- [11] X. Yang, X. Yu, H. Huang, and H. Zhu, "Energy efficiency based joint computation offloading and resource allocation in multi-access MEC systems," *IEEE Access*, vol. 7, pp. 117 054–117 062, 2019.
- [12] J. Zhang, X. Hu, Z. Ning, E. C.-H. Ngai, L. Zhou, J. Wei, J. Cheng, and B. Hu, "Energy-latency tradeoff for energy-aware offloading in mobile edge computing networks," *IEEE Internet Things J.*, vol. 5, no. 4, pp. 2633–2645, 2017.
- [13] Y. Wu, Y. Wang, F. Zhou, and R. Q. Hu, "Computation Efficiency Maximization in OFDMA-Based Mobile Edge Computing Networks," *IEEE Commun. Lett.*, vol. 24, no. 1, pp. 159–163, 2020.
- [14] S. Bi, J. Ying, and Zhang, "Computation Rate Maximization for Wireless Powered Mobile-Edge Computing with Binary Computation Offloading," *IEEE Trans. Wireless Commun.*, vol. 17, no. 6, pp. 4177–4190, 2018.
- [15] P. Richtrik and M. Tak, "Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function," *Math. Program.*, vol. 144, no. 1-2, pp. 1–38, 2014.
- [16] M. Zhao, J. Y. Ryu, J. Lee, T. Q. Quek, and S. Feng, "Exploiting trust degree for multiple-antenna user cooperation," *IEEE Trans. Wireless Commun.*, vol. 16, no. 8, pp. 4908–4923, 2017.
- [17] M. Zhao, J. Yu, W. Li, D. Liu, S. Yao, W. Feng, C. She, and T. Q. S. Quek, "Energy-Aware Offloading in Time-Sensitive Networks with Mobile Edge Computing," *ArXiv*, vol. abs/2003.12719, 2020.