# Power-Delay Tradeoff in Multi-User Mobile-Edge Computing Systems

Yuyi Mao[†], Jun Zhang[†], S.H. Song[†], and K. B. Letaief[†*], *Fellow, IEEE*

[†]Dept. of ECE, The Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong
[*]Hamad bin Khalifa University, Doha, Qatar
Email: {ymaoac, eejzhang, eeshsong, eekhaled}@ust.hk

*Abstract*—Mobile-edge computing (MEC) has recently emerged as a promising paradigm to liberate mobile devices from increasingly intensive computation workloads, as well as to improve the quality of computation experience. In this paper, we investigate the tradeoff between two critical but conflicting objectives in multi-user MEC systems, namely, the power consumption of mobile devices and the execution delay of computation tasks. A power consumption minimization problem with task buffer stability constraints is formulated to investigate the tradeoff, and an online algorithm that decides the local execution and computation offloading policy is developed based on Lyapunov optimization. Specifically, at each time slot, the optimal frequencies of the local CPUs are obtained in closed forms, while the optimal transmit power and bandwidth allocation for computation offloading are determined with the Gauss-Seidel method. Performance analysis is conducted for the proposed algorithm, which indicates that the power consumption and execution delay obeys an $[O(1/V), O(V)]$ tradeoff with $V$ as a control parameter. Simulation results are provided to validate the theoretical analysis and demonstrate the impacts of various parameters to the system performance.

*Index Terms*—Mobile-edge computing, dynamic voltage and frequency scaling, power control, bandwidth allocation, Lyapunov optimization, quality of computation experience.

## I. Introduction

The increasing popularity of smart mobile devices is driving the development of mobile applications, which can be computation-intensive, e.g., interactive online gaming, face recognition and 3D modeling. This poses more stringent requirements on the quality of computation experience, which cannot be easily satisfied by the limited processing capability of mobile devices. As a result, new solutions to handle the explosive computation demands and the ever-increasing computation quality requirements are emerging [1]. Mobile-edge computing (MEC) is such a promising technique to release the tension between the computation-intensive applications and the resource-limited mobile devices [2]. Different from conventional cloud computing systems, where remote public clouds are utilized, MEC offers computation capability within the radio access network. Therefore, by offloading the computation tasks from the mobile devices to the MEC servers, the quality of computation experience, including the device energy consumption and execution latency, can be greatly improved [3].

Nevertheless, the efficiency of computation offloading highly depends on the wireless channel conditions, as offloading tasks requires effective data transmission. Therefore, computation offloading policies for MEC systems have attracted significant attention in recent years [4]-[11]. For applications with strict deadline requirements, the local execution energy consumption was minimized by adopting dynamic voltage and frequency scaling (DVFS) techniques, and the energy consumption for computation offloading was optimized using data transmission scheduling in [4]. In [5], joint allocation of communication and computational resources for femto-cloud computing systems was proposed, where each computation task should be completed before its deadline. In [6], a dynamic computation offloading policy was developed for MEC systems with energy harvesting devices under a strict execution delay requirement. Besides, a decentralized computation offloading algorithm was proposed to minimize the computation overhead for multi-user MEC systems in [7].

Imposing strict execution delay constraints makes the computation offloading design more tractable, as only short-term performance, e.g., the performance for executing a single task, needs to be considered. However, it may be impractical for applications that can tolerate a certain period of execution latency, such as multi-media streaming. For such type of applications, the long-term system performance is more relevant, where the coupling among the randomly arrived tasks cannot be ignored. In order to minimize the long-term average energy consumption, a stochastic control algorithm was proposed in [8], which determines the offloaded software components. In [9], a delay-optimal stochastic task scheduling algorithm was developed for single-user MEC systems. Moreover, an online task scheduling algorithm was proposed to investigate the energy-delay tradeoff for MEC systems with a multi-core mobile device in [10], and this study was later extended to scenarios with heterogeneous types of mobile applications in [11]. Unfortunately, existing works only focused on single-user MEC systems, and the design methodologies for multi-user MEC systems remain unknown.

In this paper, we consider a general MEC system with multiple mobile devices, where computation tasks arrive at the mobile devices in a stochastic manner. Joint design of local execution and computation offloading strategies will be investigated. With multiple devices, the design becomes much more challenging as intelligent management of the radio resources for computation offloading, e.g., the transmit power

and available spectrum, is needed. We formulate a power consumption minimization problem with task buffer stability constraints. An online algorithm is proposed based on Lyapunov optimization, which decides the CPU-cycle frequencies for local execution, and the transmit power and bandwidth allocation for computation offloading. In particular, the optimal CPU-cycle frequencies are obtained in closed forms, while the optimal transmit power and bandwidth allocation are determined by the Gauss-Seidel method. Performance analysis is conducted for the proposed algorithm, which explicitly characterizes the tradeoff between the power consumption of the mobile devices and the execution delay. Simulation results verify the theoretical analysis and demonstrate that the proposed algorithm is capable of controlling the power consumption and execution delay performance in multi-user MEC systems.

The organization of this paper is as follows. We introduce the system model in Section II. The power consumption minimization problem is formulated in Section III, and an online local execution and computation offloading policy is developed in Section IV. Simulation results will be shown in Section V, and we will conclude this paper in Section VI.
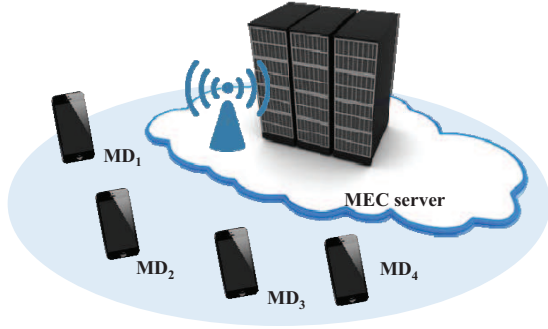
## II. SYSTEM MODEL



Fig. 1. A mobile-edge computing system with four mobile devices (MDs).

We consider a mobile-edge computing (MEC) system as shown in Fig. 1, where $N$ mobile devices running computation-intensive applications are assisted by an MEC server. The MEC server could be a small data center installed at a wireless access point deployed by the telecom operator. Therefore, it can be accessed by the mobile devices through wireless channels, and will execute the computation tasks on behalf of the mobile devices [3], [4]. By offloading part of the computation tasks to the MEC server, the mobile devices could enjoy a higher level of quality of computation experience [3].

The available system bandwidth is $w$ Hz, which is shared by the mobile devices, and the noise power spectral density at the receiver of the MEC server is denoted as $N_0$. Time is slotted and the time slot length is $\tau$. For convenience, we denote the index sets of the mobile devices and the time slots as $\mathcal{N} \triangleq \{1, \cdots, N\}$ and $\mathcal{T} \triangleq \{0, 1, \cdots\}$, respectively.

### A. Computation Task and Task Queueing Models

We assume the mobile devices are running fine-grained tasks [11]: At the beginning of the $t$th time slot, $A_i(t)$ (bits) of

computation tasks arrive at the $i$th mobile device, which can be processed starting from the $(t+1)$th time slot. Without loss of generality, we assume the $A_i(t)$'s in different time slots are independent and identically distributed (i.i.d.) within $[A_{i,\min}, A_{i,\max}]$ with $\mathbb{E}[A_i(t)] = \lambda_i, i \in \mathcal{N}$.

In each time slot, part of the computation tasks of the $i$th mobile device, denoted as $D_{l,i}(t)$, will be executed at the local CPU, while $D_{r,i}(t)$ bits of the computation tasks will be offloaded to and executed by the MEC server. The arrived but not yet executed tasks will be queued in the task buffer at each mobile device, and the queue lengths of the task buffers at the beginning of the $t$th time slot are denoted as $\mathbf{Q}(t) \triangleq [Q_1(t), \cdots, Q_N(t)]$ with $\mathbf{Q}(0) = \mathbf{0}$, where $Q_i(t)$ evolves according to the following equation:

$$Q_i(t+1) = \max\{Q_i(t) - D_{\Sigma,i}(t), 0\} + A_i(t), t \in \mathcal{T}. \quad (1)$$

In (1), $D_{\Sigma,i}(t) = D_{l,i}(t) + D_{r,i}(t)$ is the amount of tasks departing from the task buffer at the $i$th device in time slot $t$.

### B. Local Execution Model

In order to process one bit of task input at the $i$th mobile device, $L_i$ CPU cycles will be needed, which depends on the types of applications and can be obtained by off-line measurements [12]. Denote the scheduled CPU-cycle frequency for the $i$th mobile device in the $t$th time slot as $f_i(t)$, which cannot exceed $f_{i,\max}$. Thus, $D_{l,i}(t)$ can be expressed as

$$D_{l,i}(t) = \tau f_i(t) L_i^{-1}. \quad (2)$$

Accordingly, the power consumption for local execution at the $i$th mobile device is given by

$$p_{l,i}(t) = \kappa f_i^3(t), \quad (3)$$

where $\kappa$ is the effective switched capacitance related to the chip architecture [13].

### C. MEC Server Execution Model

To offload the computation tasks for MEC server execution, the input bits of the tasks need to be delivered to the MEC server. For simplicity, we assume the MEC server is equipped with an $N$-core high-speed CPU so that it can execute $N$ different applications in parallel, and the processing latency at the MEC server is negligible. We leave the investigation of more general MEC servers to our future work.

The wireless channels between the mobile devices and the MEC server are i.i.d. frequency-flat block fading. Denote the small-scale fading channel power gain from the $i$th mobile device to the MEC server at the $t$th time slot as $h_i(t)$, which is assumed to have a finite mean value, i.e., $\mathbb{E}[h_i(t)] = \overline{h_i} < \infty$. Thus, the channel power gain from the $i$th mobile device to the MEC server can be represented by $H_i(t) = h_i(t) g_0 (d_0/d_i)^\theta$, where $g_0$ is the path-loss constant, $d_0$ is the reference distance, $\theta$ is the path-loss exponent, and $d_i$ is the distance from mobile device $i$ to the MEC server. Hence, the amount of offloaded tasks at the $i$th mobile device in time slot $t$ is given by

$$D_{r,i}(t) = \begin{cases} \alpha_i(t) w\tau \log_2\left(1 + \frac{H_i(t)p_{\text{tx},i}(t)}{\alpha_i(t)N_0w}\right), & \alpha_i(t) > 0 \\ 0, & \alpha_i(t) = 0, \end{cases} \quad (4)$$

where $p_{\text{tx},i}(t)$ is the transmit power with the maximum value of $p_{i,\max}$, and $\alpha_i(t)$ is the portion of bandwidth allocated to the $i$th mobile device. Denote $\boldsymbol{\alpha}(t) \triangleq [\alpha_1(t), \cdots, \alpha_N(t)]$ as the bandwidth allocation vector, which should be chosen from the feasible set $\mathcal{A}$ [14], i.e.,

$$\boldsymbol{\alpha}(t) \in \mathcal{A} \triangleq \left\{ \boldsymbol{\alpha} \in \mathbb{R}_+^N | \sum_{i \in \mathcal{N}} \alpha_i \leq 1 \right\}. \tag{5}$$

## III. PROBLEM FORMULATION

In this section, we will first introduce the performance metrics, namely, the power consumption of the mobile devices and the average queue lengths of the task buffers. A power consumption minimization problem will then be formulated to facilitate the investigation of the power-delay tradeoff.

The average power consumption of the mobile devices, including the power consumed by the local CPUs and the transmit power for computation offloading, can be expressed as

$$\overline{P} = \lim_{T \to \infty} \frac{1}{T} \mathbb{E}\left[\sum_{t=0}^{T-1} P(t)\right], \tag{6}$$

where $P(t) \triangleq \sum_{i \in \mathcal{N}} (p_{\text{tx},i}(t) + p_{l,i}(t))$.

According to the *Little's Law* [15], the execution delay is proportional to the average queue length of the task buffer. Hence, we adopt the average queue length of the task buffer as a measurement of the execution delay, which can be written as

$$\overline{Q}_i = \lim_{T \to \infty} \frac{1}{T} \mathbb{E}\left[\sum_{t=0}^{T-1} Q_i(t)\right], i \in \mathcal{N}. \tag{7}$$

Denote $\mathbf{f}(t) \triangleq [f_1(t), \cdots, f_N(t)]$ and $\mathbf{p}_{\text{tx}}(t) \triangleq [p_{\text{tx},1}(t), \cdots, p_{\text{tx},N}(t)]$. Thus, the power consumption minimization problem is formulated as

$$\mathbf{P}_1 : \min_{\mathbf{f}(t), \mathbf{p}_{\text{tx}}(t), \boldsymbol{\alpha}(t)} \overline{P}$$
$$\text{s.t.} \quad \boldsymbol{\alpha}(t) \in \mathcal{A}, t \in \mathcal{T} \tag{8}$$
$$0 \leq f_i(t) \leq f_{i,\max}, i \in \mathcal{N}, t \in \mathcal{T} \tag{9}$$
$$0 \leq p_{\text{tx},i}(t) \leq p_{i,\max}, i \in \mathcal{N}, t \in \mathcal{T} \tag{10}$$
$$\lim_{t \to \infty} \frac{\mathbb{E}\left[|Q_i(t)|\right]}{t} = 0, i \in \mathcal{N}, \tag{11}$$

where (9) and (10) are the CPU-cycle frequency constraint and the transmit power constraint, respectively. (11) requires the task buffers to be mean rate stable [16], which ensures that all the arrived computation tasks can be executed with finite delay. In general, $\mathbf{P}_1$ is a stochastic optimization problem, for which, the CPU-cycle frequency, the transmit power as well as the bandwidth allocation need to be determined for each device at each time slot. This problem is difficult to solve as the optimal decisions are temporally correlated. Also, a joint consideration on the local execution and computation offloading strategies is needed, as both of them affect the system performance. Besides, the spatial coupling of the bandwidth allocation among different mobile devices poses another challenge.

Instead of solving $\mathbf{P}_1$ directly, we consider $\mathbf{P}_2$, which is a modified version of $\mathbf{P}_1$ by replacing set $\mathcal{A}$ in (5) by set $\tilde{\mathcal{A}}$, with $\tilde{\mathcal{A}}$ defined as

$$\tilde{\mathcal{A}} = \left\{ \boldsymbol{\alpha} \in \mathbb{R}_+^N | \sum_{i \in \mathcal{N}} \alpha_i \leq 1, \alpha_i \geq \epsilon_A, i \in \mathcal{N} \right\}. \tag{12}$$

With such modification, the departure function of MEC server execution, $D_{r,i}(t)$, is continuous and differentiable with respect to $\boldsymbol{\alpha}(t) \in \tilde{\mathcal{A}}$. In addition, the optimal value of $\mathbf{P}_2$ is larger but can be made arbitrarily close to that of $\mathbf{P}_1$ by setting $\epsilon_A$ ($\epsilon_A \in (0, 1/N)$) to be sufficiently small. Furthermore, any feasible solution for $\mathbf{P}_2$ is also feasible for $\mathbf{P}_1$. Thus, we will focus on $\mathbf{P}_2$ in the remainder of this paper.

## IV. ONLINE LOCAL EXECUTION AND COMPUTATION OFFLOADING POLICY

In this section, we will propose an online local execution and computation offloading policy to solve $\mathbf{P}_2$ based on Lyapunov optimization [16], where a deterministic problem needs to be solved at each time slot. We will then analyze the performance of the proposed algorithm and reveal the power-delay tradeoff in multi-user MEC systems.

### A. Lyapunov Optimization-Based Online Algorithm

To present the algorithm, we first define the Lyapunov function as

$$L(\mathbf{Q}(t)) = \frac{1}{2} \sum_{i \in \mathcal{N}} Q_i^2(t). \tag{13}$$

Thus, the Lyapunov drift function can be written as

$$\Delta(\mathbf{Q}(t)) = \mathbb{E}\left[L(\mathbf{Q}(t+1)) - L(\mathbf{Q}(t)) | \mathbf{Q}(t)\right]. \tag{14}$$

Accordingly, the Lyapunov drift-plus-penalty function can be expressed as

$$\Delta_V(\mathbf{Q}(t)) = \Delta(\mathbf{Q}(t)) + V \cdot \mathbb{E}\left[P(t) | \mathbf{Q}(t)\right], \tag{15}$$

where $V$ (bits$^2 \cdot \text{W}^{-1}$) is a control parameter in the proposed algorithm. We find an upper bound of $\Delta_V(\mathbf{Q}(t))$ under any feasible $\mathbf{f}(t)$, $\mathbf{p}_{\text{tx}}(t)$, and $\boldsymbol{\alpha}(t)$, as specified in Lemma 1.

*Lemma 1:* For arbitrary $\mathbf{f}(t), \mathbf{p}_{\text{tx}}(t), \boldsymbol{\alpha}(t)$ such that $\forall i \in \mathcal{N}$, $f_i(t) \in [0, f_{i,\max}]$, $p_{\text{tx},i}(t) \in [0, p_{i,\max}]$, and $\boldsymbol{\alpha}(t) \in \tilde{\mathcal{A}}$, $\Delta_V(t)$ is upper bounded by

$$\Delta_V(\mathbf{Q}(t)) \leq -\mathbb{E}\left[\sum_{i \in \mathcal{N}} Q_i(t) (D_{\Sigma,i}(t) - A_i(t)) | \mathbf{Q}(t)\right]$$
$$+ V \cdot \mathbb{E}\left[P(t) | \mathbf{Q}(t)\right] + C, \tag{16}$$

where $C$ is a constant.

*Proof:* Proof is omitted due to space limitation. ∎

The main idea of the proposed online local execution and computation offloading policy is to minimize the upper bound of $\Delta_V(\mathbf{Q}(t))$ in the right-hand side of (16) greedily at each time slot. By doing so, the amount of tasks waiting in the task buffers can be maintained at a small level. Meanwhile, the power consumption of the mobile devices can be minimized. The proposed algorithm is summarized in Algorithm 1, where a deterministic optimization problem $\mathbf{P}_{\text{PTS}}$ needs to be solved

at each time slot. It is worthy to note that the objective function of $\mathbf{P}_{\text{PTS}}$ corresponds to the right-hand side of (16), and all the constraints in $\mathbf{P}_2$ except the stability constraints in (11) are retained in $\mathbf{P}_{\text{PTS}}$. The optimal solution for $\mathbf{P}_{\text{PTS}}$ will be developed in the next subsection.

---

**Algorithm 1** Lyapunov Optimization-Based Online Local Execution and Computation Offloading Policy

---

1: At the beginning of the $t$th time slot, obtain $\{Q_i(t)\}$, $\{H_i(t)\}$, and $\{A_i(t)\}$.
2: Determine $\mathbf{f}(t)$, $\mathbf{p}_{\text{tx}}(t)$ and $\boldsymbol{\alpha}(t)$ by solving

$$\mathbf{P}_{\text{PTS}}: \min_{\mathbf{f}(t),\mathbf{p}_{\text{tx}}(t),\boldsymbol{\alpha}(t)} -\sum_{i\in\mathcal{N}} Q_i(t) D_{\Sigma,i}(t) + V \cdot P(t)$$

$$\text{s.t.} \quad \boldsymbol{\alpha}(t) \in \tilde{\mathcal{A}}, (9) \text{ and } (10).$$

3: Update $\{Q_i(t)\}$ according to (1) and set $t = t+1$.

---

### B. Optimal Solution For $\mathbf{P}_{\text{PTS}}$

In this subsection, we will derive the optimal CPU-cycle frequencies, transmit powers and bandwidth allocation vector for $\mathbf{P}_{\text{PTS}}$.

**Optimal CPU-cycle Frequencies:** It is straightforward to show that the optimal CPU-cycle frequency for the $i$th mobile device in time slot $t$ can be obtained by solving

$$\mathbf{SP}_1: \min_{0\leq f_i(t)\leq f_{i,\max}} -Q_i(t)\tau f_i(t) L_i^{-1} + V\kappa f_i^3(t), \quad (17)$$

and its optimal solution is achieved at either the stationary point of the objective function or one of the boundary points, which is given by

$$f_i^\star(t) = \min\left\{ f_{i,\max}, \sqrt{\frac{Q_i(t)\tau}{3\kappa V L_i}} \right\}, i\in\mathcal{N}. \quad (18)$$

**Remark 1:** Note that $f_i^\star(t)$ increases with $Q_i(t)$ as it is desirable to execute more tasks in order to keep the queue length of the task buffer small. Besides, $f_i^\star(t)$ decreases with both $V$ and $L_i$: With a larger value of $V$, the weight of the power consumption becomes larger, and thus the local CPU slows down its frequency to reduce power consumption; with a larger value of $L_i$, local execution becomes less efficient as more CPU cycles are needed to process per bit of task input, which leads to a smaller CPU-cycle frequency.

**Optimal Transmit Power and Bandwidth Allocation:** After decoupling $\mathbf{f}(t)$ from $\mathbf{P}_{\text{PTS}}$, the optimal $\mathbf{p}_{\text{tx}}^\star(t)$ and $\boldsymbol{\alpha}^\star(t)$ can be obtained by solving

$$\mathbf{SP}_2: \min_{\boldsymbol{\alpha}(t),\mathbf{p}_{\text{tx}}(t)} -\sum_{i\in\mathcal{N}} Q_i(t) D_{r,i}(t) + V\sum_{i\in\mathcal{N}} p_{\text{tx},i}(t)$$

$$\text{s.t.} \quad 0\leq p_{\text{tx},i}(t)\leq p_{i,\max}, i\in\mathcal{N} \quad (19)$$

$$\boldsymbol{\alpha}(t) \in \tilde{\mathcal{A}}.$$

It is not difficult to identify that $\mathbf{SP}_2$ is a convex optimization problem. However, generic convex algorithms suffer from relatively high complexity as they are developed for general

convex problems and do not make use of the problem structures [17]. Motivated by this, we propose to solve $\mathbf{SP}_2$ by optimizing the transmit power and bandwidth allocation in an alternating manner, where in each iteration, the optimal transmit powers are obtained in closed forms and the optimal bandwidth allocation is determined by the *Lagrangian method*. Since $\mathbf{SP}_2$ is jointly convex with respect to $\mathbf{p}_{\text{tx}}(t)$ and $\boldsymbol{\alpha}(t)$, and its feasible region is a Cartesian product of those of $\mathbf{p}_{\text{tx}}(t)$ and $\boldsymbol{\alpha}(t)$, the alternating minimization procedure is guaranteed to converge to the global optimal solution, which is termed as the *Gauss-Seidel method* in literature [18].

**1) Optimal Transmit Power:** For a fixed bandwidth allocation vector $\boldsymbol{\alpha}(t)$, the optimal transmit power for the $i$th mobile device can be obtained by solving

$$\mathbf{P}_{\text{PWR}}: \min_{0\leq p_{\text{tx},i}(t)\leq p_{i,\max}} -Q_i(t) D_{r,i}(t) + Vp_{\text{tx},i}(t), \quad (20)$$

whose optimal solution is achieved at either the stationary point of the objective function or one of the boundary points similar to $\mathbf{SP}_1$, and it is given in closed form by

$$p_{\text{tx},i}^\star(t) =$$
$$\min\left\{ \alpha_i(t) w \max\left\{ \frac{Q_i(t)\tau}{\ln 2 \cdot V} - \frac{N_0}{H_i(t)}, 0 \right\}, p_{i,\max} \right\}, i\in\mathcal{N}. \quad (21)$$

**2) Optimal Bandwidth Allocation:** For a fixed transmit power vector $\mathbf{p}_{\text{tx}}(t)$, the optimal bandwidth allocation can be obtained by solving the following problem:

$$\mathbf{P}_{\text{BW}}: \min_{\boldsymbol{\alpha}(t)\in\tilde{\mathcal{A}}} -\sum_{i\in\mathcal{N}} Q_i(t) D_{r,i}(t), \quad (22)$$

which is more challenging as the bandwidth allocation decision is coupled among different mobile devices. Fortunately, the Lagrangian method offers an effective solution for $\mathbf{P}_{\text{BW}}$. Specifically, the partial Lagrangian can be written as

$$\mathcal{L}(\boldsymbol{\alpha}(t),\lambda) = -\sum_{i\in\mathcal{N}} Q_i(t) D_{r,i}(t) + \lambda\left(\sum_{i\in\mathcal{N}}\alpha_i(t) - 1\right), \quad (23)$$

where $\lambda \geq 0$ is the Lagrangian multiplier associated with constraint $\sum_{i\in\mathcal{N}}\alpha_i(t) \leq 1$. Based on the Karush-Kuhn-Tucker (KKT) conditions, the optimal bandwidth allocation $\boldsymbol{\alpha}^\star(t)$ and the optimal Lagrangian multiplier $\lambda^\star$ should satisfy the following equation set:

$$\begin{cases} \alpha_i^\star(t) = \max\{\epsilon_A, \mathcal{R}_i(\lambda^\star)\}, i\in\mathcal{N}, \lambda^\star > 0 \\ \sum_{i\in\mathcal{N}}\alpha_i^\star(t) = 1. \end{cases} \quad (24)$$

In (24), if $p_{\text{tx},i}^\star(t) = 0$, $\mathcal{R}_i(\lambda) \triangleq \epsilon_A$; otherwise, $\mathcal{R}_i(\lambda)$ denotes the root of $Q_i(t)\frac{dD_{r,i}(t)}{d\alpha_i(t)} = \lambda$ for $\lambda > 0$, which is positive and unique as $\frac{dD_{r,i}(t)}{d\alpha_i(t)}$ decreases with $\alpha_i(t)$. Thus, it suggests a bisection search over $[\lambda_L, \lambda_U]$ for the optimal $\lambda^\star$, where $\lambda_L = \max_{i\in\mathcal{N}} Q_i(t)\frac{dD_{r,i}(t)}{d\alpha_i(t)}|_{\alpha_i(t)=1}$, and $\lambda_U$ satisfies $\sum_{i\in\mathcal{N}}\max\{\epsilon_A, \mathcal{R}_i(\lambda_U)\} < 1$. Hence, $\mathcal{R}_i(\lambda)$ can be obtained by a bisection search over $(0, 1]$, and the searching process for the optimal $\lambda^\star$ will be terminated when $|\sum_{i\in\mathcal{N}}\max\{\epsilon_A, \mathcal{R}_i(\lambda)\} - 1| < \xi$, where $\xi$ is the accuracy

of the algorithm. Details of the Lagrangian method for $\mathbf{P}_{\mathrm{BW}}$ are summarized in Algorithm 2.

---

**Algorithm 2** Lagrangian Method for $\mathbf{P}_{\mathrm{BW}}$
---
1: Set $\xi = 10^{-7}$, $\lambda_U = \lambda_L$, $l = 0$, $I_{\max} = 200$, $\beta = 1.5$, $\epsilon_A = 10^{-4}$.
2: Set $\alpha_i(t) = \max\{\epsilon_A, \mathcal{R}_i(\lambda_U)\}, i \in \mathcal{N}$.
3: **While** $\sum_{i \in \mathcal{N}} \alpha_i(t) \geq 1$ **do**
4:     $\lambda_U = \beta \cdot \lambda_U$.
5:     Set $\alpha_i(t) = \max\{\epsilon_A, \mathcal{R}_i(\lambda_U)\}, i \in \mathcal{N}$.
6: **Endwhile**
7: **While** $|\sum_{i \in \mathcal{N}} \alpha_i(t) - 1| \geq \xi$ and $l \leq I_{\max}$ **do**
8:     $\tilde{\lambda} = \frac{1}{2}(\lambda_L + \lambda_U)$ and $l = l + 1$.
9:     Set $\alpha_i(t) = \max\{\epsilon_A, \mathcal{R}_i(\tilde{\lambda})\}, i \in \mathcal{N}$.
10:     **If** $\sum_{i \in \mathcal{N}} \alpha_i(t) > 1$ **then**
11:         $\lambda_L = \tilde{\lambda}$.
12:     **Else**
13:         $\lambda_U = \tilde{\lambda}$.
14:     **Endif**
15: **Endwhile**

---

*Remark 2:* One main benefit of the proposed online algorithm is that it does not require prior information on the computation task arrival and wireless channel fading processes, which makes it also applicable for unpredictable environments. Besides, the proposed algorithm is of low complexity, as at each time slot, the optimal CPU-cycle frequencies are obtained in closed forms, while the computation offloading policy is determined by an efficient alternating minimization algorithm. Furthermore, as will be shown in the next subsection, the achievable performance of the proposed algorithm can be analytically characterized and thus facilitates the analysis on the power-delay tradeoff in multi-user MEC systems.

*C. Performance Analysis*

In this subsection, we will provide the main theoretical result in this paper, which characterizes the upper bounds for the power consumption of the mobile devices and the average sum queue length of the task buffers. Also, the tradeoff between the power consumption and execution delay will be revealed.

*Theorem 1:* Assume that $\mathbf{P}_2$ is feasible, we have:
- The average power consumption of the mobile devices under the proposed algorithm satisfies:

$$\overline{P} \leq P_{\Sigma}^{\mathrm{opt}} + C \cdot V^{-1}, \qquad (25)$$

where $P_{\Sigma}^{\mathrm{opt}}$ is the optimal value of $\mathbf{P}_2$.
- For arbitrary $i \in \mathcal{N}$, $Q_i(t)$ is mean rate stable.
- Suppose there exist $\epsilon > 0$ and $\Psi(\epsilon)$ ($\Psi(\epsilon) > P_{\Sigma}^{\mathrm{opt}}$) that satisfy the Slater conditions [16], then the average sum queue lengths of the task buffers satisfies:

$$\sum_{i \in \mathcal{N}} \overline{Q}_i \leq \left[C + V\left(\Psi(\epsilon) - P_{\Sigma}^{\mathrm{opt}}\right)\right] \cdot \epsilon^{-1}. \qquad (26)$$

*Proof:* Proof is omitted due to space limitation. ∎

*Remark 3:* Theorem 1 shows that under the proposed online local execution and computation offloading policy, the

worst-case power consumption of the mobile devices decreases inversely proportional to $V$, while the upper bound of the execution delay increases linearly with $V$, i.e., there exists an $[O(1/V), O(V)]$ tradeoff between these two objectives. Thus, we can balance the power consumption and execution delay by adjusting $V$: For delay-sensitive types of applications, we can use a small value of $V$; while for energy-sensitive networks and delay-tolerant applications, a large value of $V$ can be adopted.

## V. SIMULATION RESULTS

In simulations, we assume $N$ mobile devices are located at an equal distance of 150 m from the MEC server. The small-scale fading channel power gains are exponentially distributed with unit mean. Besides, $\kappa = 10^{-27}$, $\tau = 1$ ms, $w = 10$ MHz, $N_0 = -174$ dBm/Hz, $g_0 = -40$ dB, $d_0 = 1$ m, $\theta = 4$, $f_{i,\max} = 1$ GHz, $p_{i,\max} = 500$ mW, $A_i(t)$ is uniformly distributed within $[0, A_{i,\max}]$, and $L_i = 737.5$ cycles/bit, $i \in \mathcal{N}$ [12]. The simulation results are averaged over 5000 time slots.
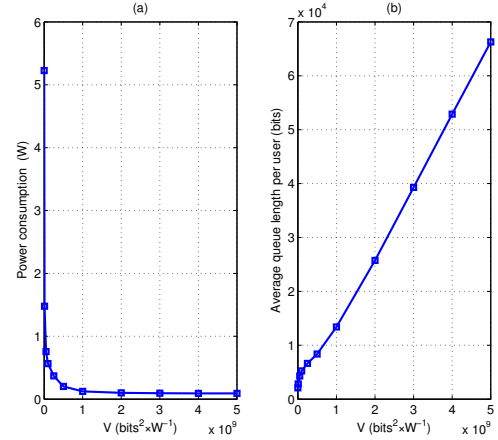


Fig. 2. Power consumption of the mobile devices/average queue length per user vs. the control parameter $V$, $N = 5$ and $A_{i,\max} = 4$ kbits.

We first show the relationship between the power consumption of the mobile devices/average queue length of the task buffers and the control parameter $V$ in Fig. 2. We see from Fig. 2a) that the power consumption decreases inversely proportional to $V$ and converges to $P_{\Sigma}^{\mathrm{opt}}$ when $V$ is sufficiently large. Meanwhile, as shown in Fig. 2b), the average queue length of the task buffers increases linearly with $V$ and becomes unbounded when $V$ goes to infinity. These results verify the $[O(1/V), O(V)]$ tradeoff between the power consumption and execution delay as shown in Theorem 1.

In Fig. 3, we show the relationship between the power consumption and execution delay for scenarios with and without MEC[1]. It is observed that by increasing $V$ from $10^6$ to $5 \times 10^9$ bit$^2 \cdot$ W$^{-1}$, the power consumption of the mobile devices decreases significantly for both cases. However, the behaviors of the execution delay are substantially different:

---
[1]The average execution delay is calculated by $\sum_{i \in \mathcal{N}} \overline{Q}_i / \sum_{i \in \mathcal{N}} \lambda_i$ (time slots) according to the Little's Law.
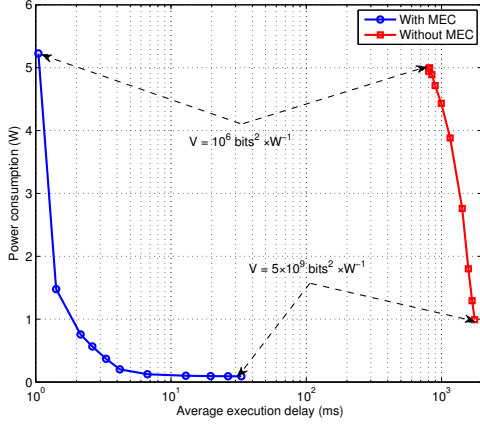
Fig. 3. Power consumption of the mobile devices vs. execution delay for systems with and without MEC, $N = 5$ and $A_{i,\max} = 4$ kbits.
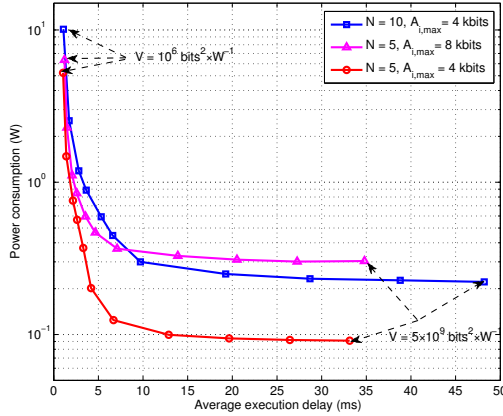


Fig. 4. Power consumption of the mobile devices vs. execution delay.

With MEC, the execution delay decreases sharply from 33.2 to 1.05 ms as $V$ decreases, while without MEC, the execution delay has minor changes at around $10^3$ ms. This is because without the aid of the MEC server, the devices cannot stabilize their task buffers even with a small $V$, where the local CPUs operate at their maximum frequencies. Therefore, we verify the benefits of MEC for improving the quality of computation experience.

By varying $A_{i,\max}$ and $N$, we show the relationship between the power consumption and execution delay in Fig. 4. In general, the average execution delay increases as the power consumption decreases, which indicates that a proper $V$ should be chosen to balance the two desirable objectives. For instance, with $N = 5$ and $A_{i,\max} = 4$ kbits, if the average execution delay requirement is 20 ms, $V = 3 \times 10^9$ bits$^2 \cdot$ W$^{-1}$ can be chosen, and the power consumption will be 0.1 W. Besides, with a given execution delay, the power consumption increases with the computation task arrival rate (the number of mobile devices), which agrees with the intuitions as the workload of the MEC system becomes heavier, more power is needed to stabilize the task buffers. In addition, when $V$ goes to infinity, doubling the computation task arrival rates results

in a higher power consumption than doubling the number of mobile devices, which is due to the increased multi-user diversity gain and the availability of extra local CPUs.

## VI. CONCLUSIONS

In this paper, we investigated the power-delay tradeoff in a multi-user mobile-edge computing system. A power consumption minimization problem with task buffer stability constraints was formulated, and an online algorithm that decides the local execution and computation offloading policy was derived based on Lyapunov optimization. Performance analysis was conducted for the proposed algorithm, which explicitly characterizes the $[O(1/V), O(V)]$ tradeoff between the power consumption and execution delay performance. Simulation results validated the theoretical analysis, and showed that the proposed algorithm is capable of balancing the power consumption of the mobile devices and the quality of computation experience. For future investigation, it would be interesting to extend the findings in this work to scenarios with fairness considerations among multiple devices.

## REFERENCES

[1] J. Gubbi, R. Buyya, S. Marusic, and M. Palaniswmi, "Internet of Things (IoT): A vision, architectural elements, and future directions," *ELSEVIER Future Gener. Comp. Syst.*, vol. 29, no. 7, pp. 1645-1660, Sep. 2013.

[2] European Telecommunications Standards Institute (ETSI), "Mobile-edge computing-Introductory technical white paper," Sep. 2014.

[3] M. Satyanarayanan, P. Bahl, R. Caceres, and N. Davies, "The case for vm-based cloudlets in mobile computing," *IEEE Pervasive Comput.*, vol. 8, no. 4, pp. 14-23, Oct. 2009.

[4] W. Zhang, Y. Wen, K. Guan, D. Kilper, H. Luo, and D. Wu, "Energy-optimal mobile cloud computing under stochastic wireless channel," *IEEE Trans. Wireless Commun.*, vol. 12, no. 9, pp. 4569-4581, Sep. 2013.

[5] O. Munoz, A. Iserte, and J. Vidal, "Optimization of radio and computational resources for energy efficiency in latency-constrained application offloading," *IEEE Trans. Veh. Technol.*, vol. 64, no. 10, pp. 4738-4755, Oct. 2015.

[6] Y. Mao, J. Zhang, and K. B. Letaief, "Dynamic computation offloading for mobile-edge computing with energy harvesting devices," *IEEE J. Sel. Areas Commun.*, to appear.

[7] X. Chen, "Decentralized computation offloading game for mobile cloud computing," *IEEE Trans. Parallel Distrib. Syst.*, vol. 26, no. 4, pp. 974-983, Apr. 2015.

[8] D. Huang, P. Wang, and D. Niyato, "A dynamic offloading algorithm for mobile computing," *IEEE Trans. Wireless Commun.*, vol. 11, no. 6, pp. 1991-1995, Jun. 2012.

[9] J. Liu, Y. Mao, J. Zhang, and K. B. Letaief, "Delay-optimal computation task scheduling for mobile-edge computing systems," in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, Barcelona, Spain, Jul. 2016.

[10] Z. Jiang and S. Mao, "Energy delay trade-off in cloud offloading for multi-core mobile devices," in *Proc. IEEE Global Commun. Conf. (GLOBECOM)*, San Diego, CA, Dec. 2015.

[11] J. Kwak, Y. Kim, J. Lee, and S. Chong, "DREAM: Dynamic resource and task allocation for energy minimization in mobile cloud systems," *IEEE J. Sel. Areas Commun.*, vol. 33, no. 12, pp. 2510-2523, Dec. 2015.

[12] A. P. Miettinen and J. K. Nurminen, "Energy efficiency of mobile clients in cloud computing," in *Proc. USENIX Conference on Hot Topics in Cloud Computing (HotCloud)*, Boston, MA, Jun. 2010.

[13] T. D. Burd and R. W. Brodersen, "Processor design for portable systems," *Kluwer J. VLSI Signal Process. Syst.*, vol. 13, no. 2/3, pp. 203-221, Aug. 1996.

[14] Z. Wang, V. Aggarwal, and X. Wang, "Joint energy-bandwidth allocation in multiple broadcast channels with energy harvesting," *IEEE Trans. Commun.*, vol. 63, no. 10, pp. 3842-3885, Oct. 2015.

[15] S. M. Ross, *Introduction to probability models*. Academic Press, 2014.

[16] M. J. Neely, *Stochastic network optimization with application to communication and queueing systems*. Morgan & Calypool, 2010.

[17] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge University Press, 2004.

[18] L. Grippo and M. Sciandron, "On the convergence of the block nonlinear Gauss-Seidel method under convex constraints," *ELSEVIER Oper. Res. Lett.*, vol. 26, no. 3, pp. 127-136, Apr. 2000.