

Minimising Offloading Latency for Edge-Cloud Systems with Ultra-Reliable and Low-Latency Communications

Dang Van Huynh*, Van-Dinh Nguyen[†], Saeed R. Khosravirad[‡] and Trung Q. Duong*

*Queen's University Belfast, UK (e-mail: {dhuynh01, trung.q.duong}@qub.ac.uk)

[†]University of Luxembourg, Luxembourg (e-mail: dinh.nguyen@uni.lu)

[‡]Nokia Bell Labs, USA (e-mail: saeed.khosravirad@nokia-bell-labs.com)

Abstract—We study a joint communication and computation offloading (JCCO) for hierarchical edge-cloud systems with ultra-reliable and low latency communications (URLLC). We aim to minimize the worst-case end-to-end (e2e) latency of computational tasks among multiple industrial Internet of Things (IIoT) devices by jointly optimizing offloading probabilities, processing rates, user association policies and power control subject to their service delay and energy consumption requirements as well as queueing stability conditions. To tackle the problem, we first decompose the original problem into two subproblems and then leverage the alternating optimization (AO) approach to solve them in an iterative fashion by developing newly convex approximate functions. The numerical results are provided to demonstrate the effectiveness of the proposed algorithms in terms of the e2e latency and convergence speed.

I. INTRODUCTION

Task offloading designs for mobile-edge computing (MEC) have been widely investigated in the literature (see [1] and the references therein). In most existing works, energy-efficiency and delay-efficiency are considered as major figure-of-merit in designing task offloading schemes for MEC systems [2]. In particular, an offload forwarding scheme where fog servers (FSs) cooperate with each other to tackle their heterogeneity in terms of commutation capacity and resources, improving the efficiency of power usage, was proposed in [2]. To guarantee low-latency wireless communication, short packets to convey a small amount of data must be used. This however will pose several challenges to design and optimize the performance of short packet-enabled networks since it demands for more resources (e.g., parity, redundancy) and ultrahigh reliability. Since then, resource allocation in the URLLC-based short block-length regime has recently studied to reduce the required bandwidth, the packet dropping, and maximize the energy efficiency (EE) [3], [4]. Focusing on designing URLLC-aware optimization for task offloading, the exponential-weight algorithm to balance URLLC constraints and energy consumption through online learning was investigated in [5]. The user-server association policy to reduce users' power consumption while trading off the resource allocations

for local computation and task offloading was considered in [6].

Although ESs and FSs are often equipped with more powerful computing capability than end users, they are still limited compared to large-scale cloud data centers (CDCs) at the cloud server. Against the above background, this work proposes a novel joint communication and computation for URLLC-enabled hierarchical edge-cloud system, taking into account all the above issues. Our main contributions are summarised as follows. Firstly, we formulate the worst-case end-to-end (e2e) latency minimisation problem with three layers edge-cloud systems, which takes into account various aspects of joint communication and computation offloading. Then, we proposed the OA-IA based solution to deal with the problem. Finally, extensive numerical results are provided to evaluate the effectiveness of the proposed solution compared with other scheme such as random user association (RUA), fixed power, fixed frequency allocation, and without cloud (w/o Cloud).

II. SYSTEM MODEL AND PROBLEM FORMULATION

We consider a hierarchical edge-cloud system consisting of the set $\mathcal{M} = \{1, 2, \dots, M\}$ of M UEs (IIoT devices) randomly distributed in a factory automation scenario, and the set $\mathcal{K} = \{1, 2, \dots, K\}$ of K ESs at the edge layer. Each ES is co-located with an access point (AP) to communicate with UEs over URLLC wireless links. The ESs connect with the cloud server via wired fronthaul links. **User layer** includes multiple UEs. A task of UE $m \in \mathcal{M}$ can be executed locally with processing rate f_m^{lo} (cycles/s) or offloaded to a ES with probability of $\alpha_m \in [0, 1]$. We use the indicator vector $\pi \triangleq [\pi_{mk}]_{\forall m,k}$ to denote the association between UEs and ESs. We assume that the tasks of a UE is only offloaded to one ES, i.e., $\sum_{k \in \mathcal{K}} \pi_{mk} = 1, \forall m$. **Edge layer** consists of K ESs placed close to UEs, where the processing rate of ES k is denoted as f_k^{es} (cycles/s). To minimize the processing latency while admitting computation tasks from multiple UEs, ES k can offload a portion of $\beta_{mk} \in [0, 1]$ of tasks of UE m to the cloud server through the fronthaul link. **Cloud layer** contains large-scale cloud data centers equipped with powerful processors having very high processing rate f^{cs} (cycles/s).

Suppose a task of UE m is characterised by a tuple $I_m \triangleq (D_m, C_m, T_m^{\text{max}})$, in which D_m , C_m and T_m^{max} are the input

This work was supported in part by the U.K. Royal Academy of Engineering (RAEng) under the RAEng Research Chair and Senior Research Fellowship scheme Grant RCSR2021\11\41. The work of V.-D. Nguyen was supported in part by the ERC AGNOSTIC project, ref. H2020/ERC2020POC/957570.

data size, the required computation resource (number of CPU cycles) and the maximum delay requirement, respectively. The task arrival rate of UE m is denoted as λ_m^{lo} (tasks/s).

A. Communication Model

Each AP is equipped with $L > 1$ antennas while each UE has single antenna. The channel vector between UE m and AP k , denoted by $\mathbf{h}_{mk} \in \mathbb{C}^{L \times 1}$, can be modeled as $\mathbf{h}_{mk} = \sqrt{g_{mk}} \tilde{\mathbf{h}}_{mk}$, where g_{mk} is the large-scale channel coefficient including the pathloss and shadowing which is normalized by the noise power, and $\tilde{\mathbf{h}}_{mk}$ is the small-scale fading following the Rayleigh fading model as $\tilde{\mathbf{h}}_{mk} \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_L)$. Under a shared wireless medium, the $L \times 1$ received signal vector at AP k can be expressed as $\mathbf{y}_k = \sum_{m \in \mathcal{M}} \mathbf{h}_{mk} \sqrt{p_m} s_m + \mathbf{z}_k$, where p_m and s_m are the transmit power and unit-power data symbol of UE m , respectively; $\mathbf{z}_k \sim \mathcal{CN}(\mathbf{0}, \mathbf{I}_L)$ is the additive white Gaussian background noise (AWGN).

Uplink Channel Estimation: The MMSE channel estimate of \mathbf{h}_{mk} is given by [7]:

$$\hat{\mathbf{h}}_{mk} = \frac{g_{mk} M p_m^p}{g_{mk} M p_m^p + 1} \mathbf{y}_{mk}^p \quad (1)$$

which follows the distribution of $\mathcal{CN}(\mathbf{0}, \sigma_{mk}^2 \mathbf{I})$, where σ_{mk}^2 is given as $\sigma_{mk}^2 = g_{mk}^2 M p_m^p / (g_{mk} M p_m^p + 1)$ and p_m^p is the pilot transmit power of UE m . According to the MMSE estimation property, the channel estimation error $\tilde{\mathbf{h}}_{mk} = \mathbf{h}_{mk} - \hat{\mathbf{h}}_{mk}$ is independent of $\hat{\mathbf{h}}_{mk}$ that follows the distribution of $\mathcal{CN}(\mathbf{0}, \delta_{mk}^2 \mathbf{I}_L)$, where δ_{mk}^2 is given by $\delta_{mk}^2 = g_{mk} / (g_{mk} M p_m^p + 1)$.

URLLC Uplink Transmission Rate: By employing the successive interference cancellation for signal detection in the uplink, the instantaneous signal-to-interference-plus-noise (SINR) of UE m can be expressed as

$$\gamma_m(\mathbf{p}, \boldsymbol{\pi}_m) = \frac{\sum_{k \in \mathcal{K}} \pi_{mk} p_m \|\hat{\mathbf{h}}_{mk}\|^4}{\Phi_m(\mathbf{p}, \boldsymbol{\pi}_m)} \quad (2)$$

where

$$\Phi_m(\mathbf{p}, \boldsymbol{\pi}_m) \triangleq \sum_{i > m} \sum_{k \in \mathcal{K}} \pi_{mk} p_i |\hat{\mathbf{h}}_{mk}^H \hat{\mathbf{h}}_{ik}|^2 + \sum_{i \in \mathcal{M}} \sum_{k \in \mathcal{K}} (1 - \pi_{mk}) p_i \times |\hat{\mathbf{h}}_{mk}^H \hat{\mathbf{h}}_{ik}|^2 + \sum_{i \in \mathcal{M}} \sum_{k \in \mathcal{K}} \pi_{mk} p_i |\hat{\mathbf{h}}_{mk} \tilde{\mathbf{h}}_{ik}|^2 + \|\hat{\mathbf{h}}_{mk}\|^2$$

with $\mathbf{p} = \{p_m\}_{m \in \mathcal{M}}$ and $\boldsymbol{\pi}_m = \{\pi_{mk}\}_{k \in \mathcal{K}}$. In this paper, we focus the ergodic achievable rate of UE m , where $\bar{\gamma}_m(\mathbf{p}, \boldsymbol{\pi}_m) = \mathbb{E}\{\gamma_m(\mathbf{p}, \boldsymbol{\pi}_m)\}$ can be approximated by

$$\bar{\gamma}_m(\mathbf{p}, \boldsymbol{\pi}_m) = \frac{\sum_{k \in \mathcal{K}} \pi_{mk} (L - 1) p_m \sigma_{mk}^2}{\bar{\Phi}_m(\mathbf{p}, \boldsymbol{\pi}_m)} \quad (3)$$

with $\bar{\Phi}_m(\mathbf{p}, \boldsymbol{\pi}_m) \triangleq \sum_{i > m} \sum_{k \in \mathcal{K}} \pi_{mk} p_i \sigma_{i,k}^2 + \sum_{i \in \mathcal{M} \setminus m} \sum_{k \in \mathcal{K}} (1 - \pi_{mk}) p_i \sigma_{i,k}^2 + \sum_{i \in \mathcal{M}} \sum_{k \in \mathcal{K}} \pi_{mk} p_i \delta_{i,k}^2 + 1$, whose derivation is based on [7]

The uplink achievable data rate of UE m (in bits/s) under URLLC finite blocklength can be approximated as:

$$R_m(\mathbf{p}, \boldsymbol{\pi}_m) = (1 - \omega) B \log_2 [1 + \bar{\gamma}_m(\mathbf{p}, \boldsymbol{\pi}_m)] - B \sqrt{\frac{(1 - \omega) V_m(\mathbf{p}, \boldsymbol{\pi}_m)}{N} \frac{Q^{-1}(\epsilon_m)}{\ln 2}} \quad (4)$$

where $N = \Delta_t B$ denotes the blocklength with Δ_t being the transmission time interval and $\omega \triangleq M/N$; ϵ_m is the decoding error probability, $Q^{-1}(\cdot)$ is the inverse function of $Q(x)$, and $V(\mathbf{p}, \boldsymbol{\pi}_m)$ is the channel dispersion given by $V_m(\mathbf{p}, \boldsymbol{\pi}_m) = 1 - [1 + \bar{\gamma}_m(\mathbf{p}, \boldsymbol{\pi}_m)]^{-2}$.

B. Computation Model

Local Processing Latency: UE m can partially offload with the portion α_m of its task to the ES. The latency to process the remaining task at UE m with the processing rate f_m^{lo} is

$$t_m^{\text{lo}}(\alpha_m, f_m^{\text{lo}}) = \frac{(1 - \alpha_m) C_m}{f_m^{\text{lo}}}. \quad (5)$$

Wireless transmission latency: Given the uplink data rate in (4), the latency to transmit the portion α_m of UE m 's task is calculated as

$$t_m^{\text{co}}(\alpha_m, \mathbf{p}, \boldsymbol{\pi}_m) = \frac{\alpha_m D_m}{R_m(\mathbf{p}, \boldsymbol{\pi}_m)}. \quad (6)$$

ES Processing Latency: Let λ_k^{es} and λ_m^{lo} be the arrival rates of tasks at ES k and UE m , respectively. We have $\lambda_k^{\text{es}} = \sum_{m \in \mathcal{M}} \pi_{mk} \alpha_m \lambda_m^{\text{lo}}$ [8]. We denote by $\beta_{mk} \in [0, 1]$ the offloading portion of task m from ES k to CS. As a result, the task rate to process the remaining tasks offloaded from all UEs at ES k can be computed as $\mu_k^{\text{es}} = f_k^{\text{es}} / \sum_{m \in \mathcal{M}} \pi_{mk} \alpha_m (1 - \beta_{mk}) C_m$. By following the standard queuing model M/M/1 [8], we can compute the worst-case processing latency among ESs as

$$t^{\text{es}}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) = \max_{\forall k \in \mathcal{K}} \left\{ \frac{1}{\mu_k^{\text{es}} - \lambda_k^{\text{es}}} \right\} \quad (7)$$

where $\boldsymbol{\alpha} \triangleq \{\alpha_m\}_{m \in \mathcal{M}}$ and $\boldsymbol{\beta} \triangleq \{\beta_{mk}\}_{m \in \mathcal{M}, k \in \mathcal{K}}$.

Fronthaul Transmission Latency: Each ES k transmits the portion β_{mk} of the offloaded task $\pi_{mk} \alpha_m$ of all UEs $m \in \mathcal{M}$ to CS for further processing. The worst-case transmission latency to offload tasks from ESs to CS via fronthaul links can be expressed as

$$t^{\text{fh}}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) = \max_{\forall k \in \mathcal{K}} \left\{ \sum_{m \in \mathcal{M}} \pi_{mk} \alpha_m \beta_{mk} \frac{D_m}{R_k^{\text{fh}}} \right\} \quad (8)$$

where R_k^{fh} is the fronthaul capacity between ES k and CS.

Cloud Processing Latency: The latency for the CS to process offloaded tasks ESs can be expressed as

$$t^{\text{cs}}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) = \frac{1}{\mu^{\text{cs}} - \lambda^{\text{cs}}} \quad (9)$$

where $\mu^{\text{cs}} = f^{\text{cs}} / \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}} \pi_{mk} \alpha_m \beta_{mk} C_m$ and $\lambda^{\text{cs}} = \sum_{m \in \mathcal{M}} \sum_{k \in \mathcal{K}} \pi_{mk} \alpha_m \beta_{mk} \lambda_m^{\text{lo}}$ are considered as the task rate and the task arrival rate at CS, respectively.

C. Problem Formulation

From (5)–(9), the overall e2e latency of UE m is given by

$$t_m(f_m^{\text{lo}}, \mathbf{p}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) = t_m^{\text{lo}}(\alpha_m, f_m^{\text{lo}}) + t_m^{\text{co}}(\alpha_m, \mathbf{p}, \boldsymbol{\pi}_m) + t^{\text{es}}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) + t^{\text{fh}}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}) + t^{\text{cs}}(\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}). \quad (10)$$

The total energy of UE m consumed for the local processing and uplink transmission can be computed as

$$E_m(\alpha_m, f_m^{\text{lo}}, \mathbf{p}, \boldsymbol{\pi}) = (1 - \alpha_m) \frac{\theta_m}{2} C_m(f_m^{\text{lo}})^2 + p_m \frac{\alpha_m D_m}{R_m(\mathbf{p}, \boldsymbol{\pi}_m)} \quad (11)$$

where the constant $\theta_m/2$ denotes the average switched capacitance and the average activity factor of UE m [6].

The joint communication and computation offloading (JCCO) problem that aims to minimize the worst-case e2e latency among computational tasks under their service delay and energy consumption requirements is formulated as

$$\text{JCCO: minimize } \max_{\alpha, \beta, \boldsymbol{\pi}, \mathbf{p}, \mathbf{f}} \{t_m(f_m^{\text{lo}}, \mathbf{p}, \alpha, \beta, \boldsymbol{\pi})\} \quad (12a)$$

$$\text{s.t. } t_m(f_m^{\text{lo}}, \mathbf{p}, \alpha, \beta, \boldsymbol{\pi}) \leq T_m^{\text{max}}, \forall m \quad (12b)$$

$$E_m(\alpha_m, f_m^{\text{lo}}, \mathbf{p}, \boldsymbol{\pi}) \leq E_m^{\text{max}}, \forall m \quad (12c)$$

$$R_m(\mathbf{p}, \boldsymbol{\pi}_m) \geq \sum_{k \in \mathcal{K}} \pi_{mk} R_m^{\text{min}}, \forall m \quad (12d)$$

$$\lambda_k^{\text{es}} \leq \mu_k^{\text{es}}, \forall k \quad (12e)$$

$$\lambda^{\text{cs}} \leq \mu^{\text{cs}}, \quad (12f)$$

$$\alpha, \beta \in \mathcal{D}, \boldsymbol{\pi} \in \Pi, \mathbf{p} \in \mathcal{P}, \mathbf{f} \in \mathcal{F} \quad (12g)$$

where $\mathcal{D} \triangleq \{\alpha_m, \beta_{mk}, \forall m, k | 0 \leq \alpha_m \leq 1, 0 \leq \beta_{mk} \leq 1, \forall m, k\}$, $\mathcal{P} \triangleq \{p_m, \forall m | 0 \leq p_m \leq P_m^{\text{max}}, \forall m\}$, $\mathcal{F} \triangleq \{f_m^{\text{lo}}, \forall m | 0 \leq f_m^{\text{lo}} \leq F_m^{\text{max}}, \forall m\}$, and $\Pi \triangleq \{\pi_{mk}, \forall m, k | \pi_{mk} \in \{0, 1\} \& \sum_{k \in \mathcal{K}} \pi_{mk} = 1, \forall m, k\}$ are the set constraints of offloading decisions, uplink transmission power, processing rates and association policies, respectively; Herein, P_m^{max} and F_m^{max} are the maximum power budget and processing rate of UE $m \in \mathcal{M}$, respectively. Constraints (12b) and (12c) are imposed to ensure that the overall e2e latency and energy consumption of UE m are limited by the predetermined thresholds T_m^{max} and E_m^{max} , respectively. Constraint (12d) guarantees the minimum rate requirement R_m^{min} for all UEs. Finally, constraints (12e) and (12f) are added to ensure the queue stability at ESs and CS, respectively.

III. PROPOSED AO-BASED ALGORITHMS FOR SOLVING PROBLEM JCCO

A. Approximate Convex Problems

Let us start by rewriting the JCCO problem (12) equivalently as

$$\text{minimize } \max_{\alpha, \beta, \boldsymbol{\pi}, \mathbf{p}, \mathbf{f}, \boldsymbol{\tau}} \{t_m(f_m^{\text{lo}}, \boldsymbol{\tau})\} \quad (13a)$$

$$\text{s.t. } (12c), (12d), (12e), (12f), (12g) \quad (13b)$$

$$t_m(f_m^{\text{lo}}, \boldsymbol{\tau}) \leq T_m^{\text{max}}, \forall m \quad (13c)$$

$$\tau^{\text{co}} \geq t_m^{\text{co}}(\alpha_m, \mathbf{p}, \boldsymbol{\pi}_m), \forall m \quad (13d)$$

$$\tau^{\text{es}} \geq t^{\text{es}}(\alpha, \beta, \boldsymbol{\pi}) \quad (13e)$$

$$\tau^{\text{fh}} \geq t^{\text{fh}}(\alpha, \beta, \boldsymbol{\pi}) \quad (13f)$$

$$\tau^{\text{cs}} \geq t^{\text{cs}}(\alpha, \beta, \boldsymbol{\pi}) \quad (13g)$$

where $t_m(f_m^{\text{lo}}, \boldsymbol{\tau}) \triangleq \frac{(1 - \alpha_m) C_m}{f_m^{\text{lo}}} + \tau^{\text{co}} + \tau^{\text{es}} + \tau^{\text{fh}} + \tau^{\text{cs}}$, and $\boldsymbol{\tau} \triangleq \{\tau^{\text{co}}, \tau^{\text{es}}, \tau^{\text{fh}}, \tau^{\text{cs}}\}$ are newly introduced variables

to simplify the objective function. Constraint (13c) is derived from (12b).

It is clear that the objective (13a) is a convex function in $(f_m^{\text{lo}}, \boldsymbol{\tau})$. We also note that a direct application of IA method is still inapplicable due to strong coupling between variables. Considering the fact that the decision variables (\mathbf{p}, \mathbf{f}) and $(\alpha, \beta, \boldsymbol{\pi})$ can be executed from UEs' and ESs' sides, respectively. Let us denote by $x^{(i)}$ the feasible point of x at the i -th iteration of the proposed iterative algorithm, which is a constant. By leveraging AO method, at iteration i , we decompose problem (13) into two subproblems (SPs) as follows:

$$\text{SP-1: minimize } \max_{\mathbf{p}, \mathbf{f}, \boldsymbol{\tau} | \alpha^{(i)}, \beta^{(i)}, \boldsymbol{\pi}^{(i)}} \{t_m(f_m^{\text{lo}}, \boldsymbol{\tau})\} \quad (14a)$$

$$\text{s.t. } (12c), (12d), (13c), (13d) \quad (14b)$$

$$\mathbf{p} \in \mathcal{P}, \mathbf{f} \in \mathcal{F} \quad (14c)$$

and

$$\text{SP-2: minimize } \max_{\alpha, \beta, \boldsymbol{\pi}, \boldsymbol{\tau} | \mathbf{p}^{(i+1)}, \mathbf{f}^{(i+1)}} \{t_m(f_m^{\text{lo}}, \boldsymbol{\tau})\} \quad (15a)$$

$$\text{s.t. } (12c) - (12f), (13c) - (13g) \quad (15b)$$

$$\alpha, \beta \in \mathcal{D}, \boldsymbol{\pi} \in \Pi. \quad (15c)$$

In an AO-based iterative algorithm, we first solve SP-1 for given $(\alpha^{(i)}, \beta^{(i)}, \boldsymbol{\pi}^{(i)})$ to generate the next optimal point of $(\mathbf{p}^{(i+1)}, \mathbf{f}^{(i+1)})$ and then solve SP-2 for updated value of $(\mathbf{p}^{(i+1)}, \mathbf{f}^{(i+1)})$ to generate the next feasible point $(\alpha^{(i+1)}, \beta^{(i+1)}, \boldsymbol{\pi}^{(i+1)})$. This procedure is repeated until convergence.

1) *Approximate Convex Program for SP-1*: In problem (14), non-convex parts include (12c), (12d) and (13d). Let us handle constraint (12d) first by convexifying $R_m(\mathbf{p}, \boldsymbol{\pi}_m)$ as

$$R_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) \geq \mathcal{R}_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) = \frac{(1 - \omega) B}{\ln 2} \left[\mathcal{G}_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) - \kappa_m \mathcal{V}_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) \right], \quad (16)$$

which is fully presented in the Appendix; Here $\kappa_m = Q^{-1}(\epsilon_m) / \sqrt{(1 - \omega) N}$ and $G_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) = \ln(1 + \bar{\gamma}_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}))$. As a result, constraint (12d) is iteratively replaced by the following convex constraint

$$\mathcal{R}_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) \geq \sum_{k \in \mathcal{K}} \pi_{mk}^{(i)} R_m^{\text{min}}, \forall m \quad (17)$$

Next, we introduce new variables $\mathbf{r} \triangleq \{r_m\}_{\forall m}$ to express constraint (12c) equivalently as

$$\begin{cases} (1 - \alpha_m^{(i)}) \frac{\theta_m}{2} C_m(f_m^{\text{lo}})^2 + \alpha_m^{(i)} D_m p_m r_m \leq E_m^{\text{max}}, & (18a) \\ \frac{1}{r_m} \leq \mathcal{R}_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}), \forall m & (18b) \end{cases}$$

where constraint (18a) is non-convex due to the product of $p_m r_m$. We note that $p_m r_m$ is a concave function which can be innerly approximated for $x = p_m$, $y = r_m$, $\bar{x} = p_m^{(i)}$, $\bar{y} = r_m^{(i)}$, yielding

$$(1 - \alpha_m^{(i)}) \frac{\theta_m}{2} C_m(f_m^{\text{lo}})^2 + \frac{1}{2} \alpha_m^{(i)} D_m \left(\frac{r_m^{(i)}}{p_m^{(i)}} p_m^2 + \frac{p_m^{(i)}}{r_m^{(i)}} r_m^2 \right) \leq E_m^{\text{max}}, \forall m. \quad (19)$$

Lastly, by (18b), constraint (13d) is iteratively replaced by the following linear constraint

$$\tau^{\text{co}} \geq \alpha_m^{(i)} D_m r_m, \forall m. \quad (20)$$

As a result, we obtain the following approximate convex program of SP-1 (14) solved at iteration i :

$$\text{SP-1 CVX: } \underset{\mathbf{p}, \mathbf{f}, \mathbf{r}, \boldsymbol{\tau} | \boldsymbol{\alpha}^{(i)}, \boldsymbol{\beta}^{(i)}, \boldsymbol{\pi}^{(i)}}{\text{minimize}} \quad \max_{\forall m \in \mathcal{M}} \{t_m(f_m^{\text{lo}}, \boldsymbol{\tau})\} \quad (21a)$$

$$\text{s.t. (13c), (14c), (17), (18b), (19), (20).} \quad (21b)$$

2) *Approximate Convex Program for SP-2:* For given $(\mathbf{p}^{(i+1)}, \mathbf{f}^{(i+1)})$ obtained by solving (21), we are now in position to convexify (15). To bypass the binary nature of (15), we first relax $\boldsymbol{\pi}$ to be continuous, i.e., $\boldsymbol{\pi} \in \tilde{\Pi} \triangleq \{\pi_{mk}, \forall m, k | 0 \leq \pi_{mk} \leq 1 \ \& \ \sum_{k \in \mathcal{K}} \pi_{mk} = 1, \forall m, k\}$ and rewrite it as

$$\text{SP-2: } \underset{\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\pi}, \boldsymbol{\tau} | \mathbf{p}^{(i+1)}, \mathbf{f}^{(i+1)}}{\text{minimize}} \quad \max_{\forall m \in \mathcal{M}} \{t_m(f_m^{\text{lo}, (i+1)}, \boldsymbol{\tau})\} \quad (22a)$$

$$\text{s.t. (12c) - (12f), (13c), -(13g)} \quad (22b)$$

$$\boldsymbol{\alpha}, \boldsymbol{\beta} \in \mathcal{D}, \boldsymbol{\pi} \in \tilde{\Pi}. \quad (22c)$$

Constraints (13c) and (22c) are linear while others are non-convex.

Convexity of (12d) and (12c): We rewrite $\bar{\gamma}_m(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m) = \frac{\sum_{k \in \mathcal{K}} \pi_{mk} \sigma_{mk}^2}{\bar{q}_m(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m)}$ where $\bar{q}_m(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m)$ is defined as

$$\bar{q}_m(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m) \triangleq \frac{\bar{\Phi}_m(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m)}{p_m^{(i+1)}(L-1)}.$$

By processing similarly as in SP-1, we have

$$R_m(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m) \geq \tilde{R}_m^{(i)}(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m) = \frac{(1-\omega)B}{\ln 2} \times [\tilde{\mathcal{G}}_m^{(i)}(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m) - \kappa_m \tilde{\mathcal{V}}_m^{(i)}(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m)] \quad (23)$$

where $\tilde{\mathcal{G}}_m^{(i)}(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m)$ and $\tilde{\mathcal{V}}_m^{(i)}(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m)$ are conducted similarly as in approximate convex program of SP-1. As a result, we innerly approximate constraint (12d) as

$$\tilde{R}_m^{(i)}(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m) \geq \sum_{k \in \mathcal{K}} \pi_{mk} R_m^{\min}, \forall m. \quad (24)$$

The constraint (12c) is equivalent to

$$\left\{ \begin{array}{l} \frac{1}{\tilde{R}_m^{(i)}(\mathbf{p}^{(i+1)}, \boldsymbol{\pi}_m)} \leq r_m, \\ (1-\alpha_m) \frac{\theta_m}{2} C_m(f_m^{\text{lo}})^2 + p_m D_m \alpha_m r_m \leq E_m^{\max}, \end{array} \right. \quad (25a)$$

$$(1-\alpha_m) \frac{\theta_m}{2} C_m(f_m^{(i+1)})^2 + \frac{1}{2} p_m^{(i+1)} D_m \left(\frac{r_m^{(i+1)}}{\alpha_m^{(i)}} \alpha_m^2 + \frac{\alpha_m^{(i)}}{r_m^{(i+1)}} r_m^2 \right) \leq E_m^{\max}, \forall m. \quad (25b)$$

where $\mathbf{r} \triangleq \{r_m\}_{\forall m}$ were defined (18). We can approximate $\alpha_m r_m$ in (25b) as

$$(1-\alpha_m) \frac{\theta_m}{2} C_m(f_m^{(i+1)})^2 + \frac{1}{2} p_m^{(i+1)} D_m \left(\frac{r_m^{(i+1)}}{\alpha_m^{(i)}} \alpha_m^2 + \frac{\alpha_m^{(i)}}{r_m^{(i+1)}} r_m^2 \right) \leq E_m^{\max}, \forall m. \quad (26)$$

Convexity of (12e) and (12f): By introducing new variables $\check{\phi} \triangleq \{\check{\phi}_{mk}\}_{\forall m, k}$, constraint (12e) is expressed as

$$\left\{ \begin{array}{l} \sum_{m \in \mathcal{M}} \lambda_m^{\text{lo}} \pi_{mk} \alpha_m \leq \frac{f_k^{\text{es}}}{\sum_{m \in \mathcal{M}} C_m \check{\phi}_{mk}^2}, \forall k \\ \check{\phi}_{mk}^2 \geq \pi_{mk} \alpha_m (1 - \beta_{mk}), \forall m, k. \end{array} \right. \quad (27a)$$

$$\quad (27b)$$

We iteratively replace (27a) by

$$\sum_{m \in \mathcal{M}} \frac{1}{2} \lambda_m^{\text{lo}} \left(\frac{\alpha_m^{(i)}}{\pi_{mk}^{(i)}} \pi_{mk}^2 + \frac{\pi_{mk}^{(i)}}{\alpha_m^{(i)}} \alpha_m^2 \right) \leq f_k^{\text{es}} \left(\frac{2}{\sum_{m \in \mathcal{M}} C_m (\check{\phi}_{mk}^{(i)})^2} - \frac{\sum_{m \in \mathcal{M}} C_m \check{\phi}_{mk}^2}{(\sum_{m \in \mathcal{M}} C_m (\check{\phi}_{mk}^{(i)})^2)^2} \right), \forall k \quad (28)$$

which is a convex constraint. To handle constraint (27b), we first rewrite as $\frac{\check{\phi}_{mk}^2}{1 - \beta_{mk}} \geq \pi_{mk} \alpha_m$, and approximate both sides as

$$\frac{2\check{\phi}_{mk}^{(i)} \check{\phi}_{mk}}{1 - \beta_{mk}^{(i)}} - \frac{(\check{\phi}_{mk}^{(i)})^2 (1 - \beta_{mk})}{(1 - \beta_{mk}^{(i)})^2} \geq \frac{1}{2} \left(\frac{\alpha_m^{(i)}}{\pi_{mk}^{(i)}} \pi_{mk}^2 + \frac{\pi_{mk}^{(i)}}{\alpha_m^{(i)}} \alpha_m^2 \right), \quad (29)$$

Similarly, by introducing new variables $\hat{\phi} \triangleq \{\hat{\phi}_m\}_{\forall m}$, (12f) is equivalent to

$$\left\{ \begin{array}{l} \sum_{m \in \mathcal{M}} \lambda_m^{\text{lo}} \hat{\phi}_m^2 \leq \frac{f^{\text{cs}}}{\sum_{m \in \mathcal{M}} C_m \hat{\phi}_m^2} \\ \frac{\hat{\phi}_m^2}{\alpha_m} \geq \sum_{k \in \mathcal{K}} \pi_{mk} \beta_{mk}, \forall m \end{array} \right. \quad (30a)$$

$$\quad (30b)$$

which are approximated (detailed manipulations are omitted here due to the space limit) to yield

$$\sum_{m \in \mathcal{M}} \lambda_m^{\text{lo}} \hat{\phi}_m^2 \leq f^{\text{cs}} \left(\frac{2}{\sum_{m \in \mathcal{M}} C_m (\hat{\phi}_m^{(i)})^2} - \frac{\sum_{m \in \mathcal{M}} C_m \hat{\phi}_m^2}{(\sum_{m \in \mathcal{M}} C_m (\hat{\phi}_m^{(i)})^2)^2} \right) \quad (31)$$

$$\frac{2\hat{\phi}_m^{(i)} \hat{\phi}_m}{\alpha_m^{(i)}} - \frac{(\hat{\phi}_m^{(i)})^2 \alpha_m}{(\alpha_m^{(i)})^2} \geq \sum_{k \in \mathcal{K}} \frac{1}{2} \left(\frac{\beta_{mk}^{(i)}}{\pi_{mk}^{(i)}} \pi_{mk}^2 + \frac{\pi_{mk}^{(i)}}{\beta_{mk}^{(i)}} \beta_{mk}^2 \right), \quad (32)$$

Convexity of (13d): From (25a), we rewrite (13d) as $\tau^{\text{co}} \geq D_m \alpha_m r_m$ as

$$\tau^{\text{co}} \geq \frac{1}{2} D_m \left(\frac{r_m^{(i)}}{\alpha_m^{(i)}} \alpha_m^2 + \frac{\alpha_m^{(i)}}{r_m^{(i)}} r_m^2 \right), \forall m. \quad (33)$$

Convexity of (13e): It follows from constraint (13e) that

$$\frac{f_k^{\text{es}}}{\sum_{m \in \mathcal{M}} \pi_{mk} \alpha_m (1 - \beta_{mk}) C_m} \geq \sum_{m \in \mathcal{M}} \lambda_m^{\text{lo}} \pi_{mk} \alpha_m + \frac{1}{\tau^{\text{es}}}, \forall k \quad (34)$$

which can be transformed equivalently as

$$\frac{f_k^{\text{es}}}{\sum_{m \in \mathcal{M}} C_m \check{\phi}_{mk}^2} \geq \sum_{m \in \mathcal{M}} \lambda_m^{\text{lo}} \pi_{mk} \alpha_m + \frac{1}{\tau^{\text{es}}} \quad (35)$$

by (27b). Finally, we have

$$\begin{aligned} & f_k^{\text{es}} \left(\frac{2}{\sum_{m \in \mathcal{M}} C_m (\check{\phi}_{mk}^{(i)})^2} - \frac{\sum_{m \in \mathcal{M}} C_m \check{\phi}_{mk}^2}{(\sum_{m \in \mathcal{M}} C_m (\check{\phi}_{mk}^{(i)})^2)^2} \right) \\ & \geq \sum_{m \in \mathcal{M}} \lambda_m^{\text{lo}} \frac{1}{2} \left(\frac{\alpha_m^{(i)}}{\pi_{mk}^{(i)}} \pi_{mk}^2 + \frac{\pi_{mk}^{(i)}}{\alpha_m^{(i)}} \alpha_m^2 \right) + \frac{1}{\tau^{\text{es}}}, \forall k. \end{aligned} \quad (36)$$

Convexify of (13f): We can express constraint (13f) as

$$\begin{cases} \tau^{\text{th}} \geq \sum_{m \in \mathcal{M}} \varphi_{mk}^2 \frac{D_m}{R_k^{\text{th}}}, \forall k \\ \frac{\varphi_{mk}^2}{\beta_{mk}} \geq \pi_{mk} \alpha_m, \forall m, k \end{cases} \quad (37a)$$

where $\varphi \triangleq \{\varphi_{mk}\}_{\forall m,k}$ are new variables to tackle the product of $\pi_{mk} \alpha_m \beta_{mk}$. Constraint (37b) is non-convex. Similar to (27b), we have

$$\frac{2\varphi_{mk}^{(i)} \varphi_{mk}}{\beta_{mk}^{(i)}} - \frac{\varphi_{mk}^{2(i)} \beta_{mk}}{(\beta_{mk}^{(i)})^2} \geq \frac{1}{2} \left(\frac{\alpha_m^{(i)}}{\pi_{mk}^{(i)}} \pi_{mk}^2 + \frac{\pi_{mk}^{(i)}}{\alpha_m^{(i)}} \alpha_m^2 \right). \quad (38)$$

Convexity of (13g): We can rewrite (13g) as

$$\frac{f_{cs}}{\sum_{m \in \mathcal{M}} C_m \hat{\phi}_m^2} \geq \sum_{m \in \mathcal{M}} \lambda_m^{\text{lo}} \hat{\phi}_m^2 + \frac{1}{\tau_{cs}} \quad (39)$$

by using (30b). It follows from (31) that

$$\begin{aligned} f_{cs} \left(\frac{2}{\sum_{m \in \mathcal{M}} C_m (\hat{\phi}_m^{(i)})^2} - \frac{\sum_{m \in \mathcal{M}} C_m \hat{\phi}_m^2}{(\sum_{m \in \mathcal{M}} C_m (\hat{\phi}_m^{(i)})^2)^2} \right) \\ \geq \sum_{m \in \mathcal{M}} \lambda_m^{\text{lo}} \hat{\phi}_m^2 + \frac{1}{\tau_{cs}}. \end{aligned} \quad (40)$$

Summing up, we obtain the following approximate convex program of SP-2 solved at iteration i :

$$\text{SP-2: CVX} \quad \underset{\alpha, \beta, \pi, \tau, \phi, \hat{\phi}, \mathbf{r}, \varphi | \mathbf{p}^{(i+1)}, \mathbf{f}^{(i+1)}}{\text{minimize}} \quad \max_{\forall m \in \mathcal{M}} \{t_m(f_m^{\text{lo},(i+1)}, \tau)\} \quad (41a)$$

s.t. (13c), (22c), (24), (25a), (26), (28),

$$(29), (31), (32), (33), (36), (37a), (38), (40). \quad (41b)$$

which requires the per-iteration complexity of $\mathcal{O}(\sqrt{9M} + 4MK + 3K(3MK + 3M + 4)^2)$.

B. Proposed AO-IA based Algorithms

Let denote by $\mathcal{S}_1^{(i)} \triangleq (\mathbf{p}^{(i)}, \mathbf{f}^{(i)}, \mathbf{r}^{(i)})$ and $\mathcal{S}_2^{(i)} \triangleq (\alpha^{(i)}, \beta^{(i)}, \pi^{(i)}, \check{\phi}^{(i)}, \hat{\phi}^{(i)}, \mathbf{r}^{(i)}, \varphi^{(i)})$ the feasible sets of (21) and (41) at iteration i , respectively. The overall algorithm for solving (13) is summarized in Algorithm 1.

Algorithm 1 Proposed AO-IA based Algorithm for Solving the JCCO Problem (13)

Initialization: Set $i = 0$ and randomly generate initial feasible points $\mathcal{S}_1^{(0)}$ and $\mathcal{S}_2^{(i)}$ to constraints in (21) and (41), respectively. Set the tolerance $\varepsilon = 10^{-3}$ and the maximum number of iterations I^{max} .

1: **repeat**

2: Solve problem (21) for given $\mathcal{S}_2^{(i)}$ to obtain the optimal solution denoted by $(\mathbf{p}^*, \mathbf{f}^*, \mathbf{r}^*, \tau^*)$ and update $\mathcal{S}_1^{(i+1)} := (\mathbf{p}^*, \mathbf{f}^*, \mathbf{r}^*)$;

3: Solve problem (41) for given $\mathcal{S}_1^{(i+1)}$ to obtain the optimal solution denoted by $(\alpha^*, \beta^*, \pi^*, \check{\phi}^*, \hat{\phi}^*, \mathbf{r}^*, \varphi^*, \tau^*)$ and update $\mathcal{S}_2^{(i+1)} := (\alpha^*, \beta^*, \pi^*, \check{\phi}^*, \hat{\phi}^*, \mathbf{r}^*, \varphi^*)$;

4: Set $i := i + 1$;

5: **until** Convergence or $i > I^{\text{max}}$

6: Recover binary values of π^* : $\pi_{mk}^* = \lfloor \pi_{mk}^{(i)} + 0.5 \rfloor, \forall m, k$;

7: Repeat Steps 1-5 with fixed π^* to refine the optimal solution;

8: **Output:** $(\alpha^*, \beta^*, \pi^*, \mathbf{p}^*, \mathbf{f}^*)$.

The main drawback of solving problem (22) is that the exact binary solution of π is not guaranteed at optimum, resulting in an infeasible solution to the original problem (12). To overcome this issue, we consider Step 6 in Algorithm 1 using ceiling function to recover binary value of π as $\pi_{mk}^* = \lfloor \pi_{mk}^{(i)} + 0.5 \rfloor, \forall m, k$. In Step 7, we repeat Steps 1-5 for given π^* to minimize the performance loss due to Step 6.

IV. NUMERICAL RESULTS

We consider a small-cell scenario where all APs (ESs) and UEs are located within an area of 100×100 m [6]. ESs are located at (50, 33) and (50, 66) for $K = 2$ and (50, 20), (50, 40), (50, 60), (50, 80) for $K = 4$. The large-scale fading of the channel between UE m and AP k is modeled as $g_{mk} = 10^{\text{PL}(d_{mk})/10}$, where $\text{PL}(d_{mk}) = -35.3 - 37.6 \log_{10} d_{mk}$ denotes the path loss in dB which is a function of the distance d_{mk} [9]. The number of antennas at each AP is set to $L = 8$. We assume that all UEs have the same power budget, i.e., $P_m^{\text{max}} = 23$ dBm $\forall m$. The URLLC decoding error probability is set to $\epsilon_m = 10^{-9}, \forall m$.

We set the CPU cycles of ESs and CS to 25 and 30 Giga cycles/s, respectively. For UE m , the input data size and the required computation resource are set to $D_m = 100$ KB and $C_m = 800 \times 10^6$ (cycles), respectively. The total e2e latency requirement of each UE is given as $T_m^{\text{max}} = 2$ s $\forall m$. The arrival rate of tasks is set to $\lambda_m^{\text{lo}} = 10$ (task/s), $\forall m$ [8]. Unless specifically stated otherwise, other parameters are given as follows: System bandwidth: 10 MHz; Noise power spectral density: -174 dBm/Hz; Maximum blocklength: 100; Fronthaul capacity: 1 Gbps; Maximum energy consumption: 1 Joule; Effective capacitance coefficient: 10^{-27} Watt.s³/cycle³.

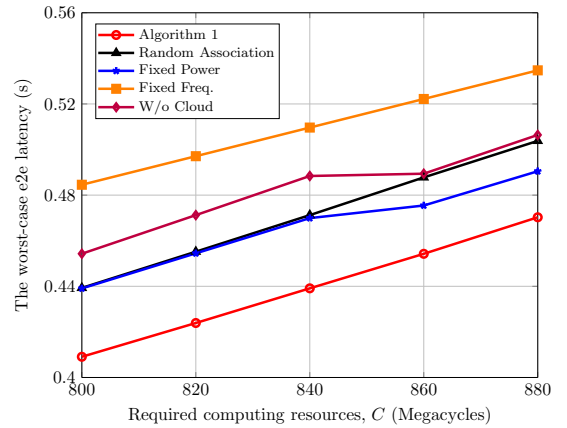


Fig. 1: The worst-case e2e latency of different resource allocation schemes versus the required computation resource $C \equiv C_m, \forall m$, with $M = 10$ UEs and $K = 2$ ESs.

Impact of required computation resource: Fig. 1 illustrates the e2e latency versus the required computation resource $C_m, \forall m$. Unsurprisingly, the e2e latency of all the considered schemes increases when C_m increases. We recall that a higher value of C_m will not only increase processing latency at UEs,

ESs and CS but also force UEs to scale down their processing rate and transmit power to satisfy constraint (12c), resulting in higher total latency. Again, Algorithm 1 still provides the lowest e2e latency amongst all the considered schemes.

Impact of offloading factor and processing rate: In hierarchical edge-cloud systems, the processing rates of ESs and CS have a strong impact on the system performance. In particular, the higher the processing rate, the higher the offloading portion from UEs to ESs. To verify this, we first show the e2e latency and average offloading portion from UEs to ESs versus ESs' processing rate $f_k^{\text{es}} \equiv f_k^{\text{es}}, \forall k$ in Fig. 2. The offloading portion of UEs slightly increases when ESs have larger computation resource, f_k^{es} . The e2e latency of Algorithm 1 sharply decreases by approximately 200 ms when the processing rate of ESs increases to 30 gigacycles/s. An important observation is that Algorithm 1 always offloads the higher portion of computation tasks, compared to sub-optimal schemes, thanks to optimal user association policies, leading to lower e2e latency.

V. CONCLUSION

We investigated the joint communication and computation task offloading in URLLC-based hierarchical edge-cloud systems. To address the practical issues of minimizing the worst-case e2e latency, we proposed an alternating optimization framework to efficiently solve the formulated problem in an iterative manner. Finally, we provided extensive numerical results to demonstrate the significant performance gain achieved by joint optimization of the communication and computation variables in hierarchical edge-cloud systems.

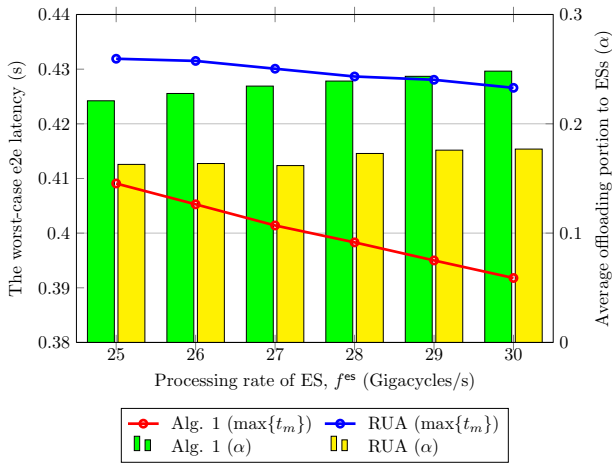


Fig. 2: The worst-case e2e latency and average offloading portion from UEs to ESs versus ESs' processing rate $f_k^{\text{es}} \equiv f_k^{\text{es}}, \forall k$, with $M = 10$ UEs and $K = 2$ ESs.

APPENDIX

We first rewrite the SINR of UE m as $\gamma_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) = p_m/q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})$. By applying [9, Eq. (72)] for $x = p_m$, $y = q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})$, $\bar{x} = p_m^{(i)}$, and $\bar{y} = q_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)})$, we have

$$G_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) \geq a_m^{(i)} - \frac{b_m^{(i)}}{p_m} - c_m^{(i)} q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) \triangleq \mathcal{G}_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})$$

$$\text{where } a_m^{(i)} = \ln\left(1 + \frac{p_m^{(i)}}{q_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})}\right) + 2\frac{p_m^{(i)}}{p_m^{(i)} + q_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})}, b_m^{(i)} = \frac{(p_m^{(i)})^2}{p_m^{(i)} + q_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})} \text{ and } c_m^{(i)} = \frac{p_m^{(i)}}{(q_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) + p_m^{(i)})q_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})}.$$

To find an upper bounding convex function approximation of $V_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})$, we apply [9, Eq. (75)] for $x = 1 - 1/(1 + \bar{\gamma}_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}))^2$ and $\bar{x} = 1 - 1/(1 + \bar{\gamma}_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)}))^2$, yielding

$$V_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) \leq d_m^{(i)} - e_m^{(i)} \frac{q_m^2(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})}{(q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) + p_m)^2} \quad (42)$$

where

$$d_m^{(i)} = 0.5\sqrt{V_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)})} + 0.5/\sqrt{V_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)})} \quad (43)$$

$$e_m^{(i)} = 0.5/\sqrt{V_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)})}. \quad (44)$$

The function $\frac{q_m^2(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})}{(q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) + p_m)^2}$ in (42) is still not convex [9], which can be further approximated as

$$\begin{aligned} \frac{q_m^2(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})}{q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) + p_m} \frac{1}{q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) + p_m} &\geq \frac{2}{q_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)}) + p_m^{(i)}} \\ &\left(\frac{2q_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)})q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})}{q_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)}) + p_m^{(i)}} - \frac{q_m^2(\mathbf{p}^{(i)})}{(q_m(\mathbf{p}^{(i)}) + p_m^{(i)})^2} \right. \\ &\left. (q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) + p_m) \right) - \frac{q_m^2(\mathbf{p}, \boldsymbol{\pi}_m^{(i)})}{(q_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) + p_m^{(i)})^2} \end{aligned}$$

By substituting this result to (42), yields

$$\begin{aligned} V_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) &\leq \mathcal{V}_m^{(i)}(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) \triangleq d_m^{(i)} - \frac{2e_m^{(i)}}{q_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)}) + p_m^{(i)}} \\ &\times \left(2f_m^{(i)}q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) - (f_m^{(i)})^2(q_m(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) + p_m) \right) \\ &+ \frac{(f_m^{(i)})^2}{q_m^2(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)})} q_m^2(\mathbf{p}, \boldsymbol{\pi}_m^{(i)}) \text{ where } f_m^{(i)} \triangleq \frac{q_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)})}{q_m(\mathbf{p}^{(i)}, \boldsymbol{\pi}_m^{(i)}) + p_m^{(i)}}. \end{aligned}$$

REFERENCES

- [1] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," *IEEE Commun. Surv. Tutorials*, vol. 19, no. 4, pp. 2322–2358, 2017.
- [2] Y. Xiao and M. Krunz, "Distributed optimization for energy-efficient fog computing in the tactile Internet," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 11, pp. 2390–2400, Nov. 2018.
- [3] C. She, C. Yang, and T. Q. S. Quek, "Joint uplink and downlink resource configuration for ultra-reliable and low-latency communications," *IEEE Trans. Commun.*, vol. 66, no. 5, pp. 2266–2280, 2018.
- [4] C. Sun, *et al.*, "Optimizing resource allocation in the short blocklength regime for ultra-reliable and low-latency communications," *IEEE Trans. Wireless Commun.*, vol. 18, no. 1, pp. 402–415, 2019.
- [5] Z. Zhou, Z. Wang, H. Yu, H. Liao, S. Mumtaz, L. Oliveira, and V. Frascolla, "Learning-based URLLC-aware task offloading for Internet of health things," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 2, pp. 396–410, Feb. 2021.
- [6] C.-F. Liu, *et al.*, "Dynamic task offloading and resource allocation for ultra-reliable low-latency edge computing," *IEEE Trans. Commun.*, vol. 67, no. 6, pp. 4132–4150, 2019.
- [7] H. Ren *et al.*, "Joint pilot and payload power allocation for massive-MIMO-enabled URLLC IIoT networks," *IEEE J. Sel. Areas Commun.*, vol. 38, no. 5, pp. 816–830, 2020.
- [8] J. Tang, W. P. Tay, T. Q. S. Quek, and B. Liang, "System cost minimization in cloud ran with limited fronthaul capacity," *IEEE Trans. Wireless Commun.*, vol. 16, no. 5, pp. 3371–3384, May 2017.
- [9] A. A. Nasir *et al.*, "Resource allocation and beamforming design in the short blocklength regime for URLLC," *IEEE Trans. Wireless Commun.*, vol. 20, no. 2, pp. 1321–1335, 2021.