

# Communication-and-Computing Latency Minimization for UAV-Enabled Virtual Reality Delivery Systems

Yi Zhou<sup>ID</sup>, Cunhua Pan<sup>ID</sup>, Member, IEEE, Phee Lep Yeoh<sup>ID</sup>, Member, IEEE,  
Kezhi Wang<sup>ID</sup>, Senior Member, IEEE, Maged Elkashlan<sup>ID</sup>, Senior Member, IEEE,  
Branka Vucetic<sup>ID</sup>, Life Fellow, IEEE, and Yonghui Li<sup>ID</sup>, Fellow, IEEE

**Abstract**—In this paper, we propose a low-latency virtual reality (VR) delivery system where an unmanned aerial vehicle (UAV) base station (U-BS) is deployed to deliver VR content from a cloud server to multiple ground VR users. Each VR input data requested by the VR users can be either projected at the U-BS before transmission or processed locally at each user. Popular VR input data is cached at the U-BS to further reduce backhaul latency from the cloud server. For this system, we design a low-complexity iterative algorithm to minimize the maximum communications and computing latency among all VR users subject to the computing, caching and transmit power constraints, which is guaranteed to converge. Numerical results indicate that our proposed algorithm can achieve a lower latency compared to other benchmark schemes. Moreover, we observe that the maximum latency mainly comes from communication latency when the bandwidth resource is limited, while it is dominated by computing latency when computing capacity is low. In addition, we find that caching is helpful to reduce latency.

**Index Terms**—UAV communication, computing, caching, latency minimization, joint optimization.

## I. INTRODUCTION

THE demand for virtual reality (VR) applications that can create high-definition ultra-immersive VR environments for mobile users has increased significantly in 5G and beyond wireless networks [1], [2]. However, due to the low computing capability and finite battery lifetime of VR users, it is extremely challenging for wireless networks to support

Manuscript received February 7, 2020; revised May 27, 2020, August 27, 2020, and October 19, 2020; accepted November 14, 2020. Date of publication November 24, 2020; date of current version March 17, 2021. The work of P. L. Yeoh was supported in part by ARC under Grant DP190100770. The work of Y. Li was supported by ARC under Grant DP190101988 and DP210103410. The work of B. Vucetic was supported in part by ARC Laureate Fellowship under Grant FL160100032. The associate editor coordinating the review of this article and approving it for publication was B. Shim. (*Corresponding author: Yonghui Li*.)

Yi Zhou, Phee Lep Yeoh, Branka Vucetic, and Yonghui Li are with the School of Electrical and Information Engineering, The University of Sydney, Sydney, NSW 2006, Australia (e-mail: yi.zhou@sydney.edu.au; phee.yeoh@sydney.edu.au; branka.vucetic@sydney.edu.au; yonghui.li@sydney.edu.au).

Cunhua Pan and Maged Elkashlan are with the School of Electronic Engineering and Computer Science, Queen Mary University of London, London E1 4NS, U.K. (e-mail: c.pan@qmul.ac.uk; maged.elkashlan@qmul.ac.uk).

Kezhi Wang is with the Department of Computer and Information Sciences, Northumbria University, Newcastle NE2 1XE, U.K. (e-mail: kezhi.wang@northumbria.ac.uk).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TCOMM.2020.3040283>.

Digital Object Identifier 10.1109/TCOMM.2020.3040283

these computing-intensive and latency-sensitive VR applications. To alleviate computing resource constraints and reduce latency, mobile edge computing (MEC) has emerged as a promising enabler for VR delivery by equipping high-capacity computing resources at the network edge [3], [4]. In [5], a task scheduling strategy was proposed to solve a transmission data consumption minimization problem with delay constraint in MEC-enabled VR systems. In [6], by optimizing the bandwidth allocation of the uplink and downlink channels, the authors solved an end-to-end delay minimization problem in a VR mobile social edge network. In [7], the authors proposed an efficient algorithm to minimize the offloading energy consumption under latency and power constraints for augmented reality applications. In [8], an energy consumption minimization framework was developed for a two-tier computing offloading MEC network. In [9], the authors implemented and developed a low-latency management framework for distributed service function chains enabling tactile internet with MEC. In [10], by jointly coordinating the task assignment, computing, and transmission resources among edge devices, multi-layer MEC servers and cloud center, the authors proposed an efficient algorithm that aimed at minimizing the system latency including total computing and transmission time in heterogeneous multi-layer MEC networks. In [11], by jointly optimizing the users' transmit power, computing capacity allocation, and user association, a latency minimization problem of an MEC system was formulated.

To further reduce latency consumption, caching has been considered for MEC servers to pre-cache popular data files from the cloud servers during off-peak periods [12]–[14]. By doing so, the backhaul latency for requesting data from the cloud server can be minimized during peak periods. In [12], the authors formulated a joint radio communication, caching and computing decision problem to maximize the average tolerant delay with a given transmission rate constraint in a fog radio access network. In [13], joint caching and computing optimization was proposed to minimize the average transmission rate in MEC-based VR delivery systems. The authors in [14] jointly optimized the computation offloading, content caching and resource allocation such that the total latency consumption is minimized.

Due to its mobility and flexibility, unmanned aerial vehicle (UAV) is an ideal platform to provide high-quality

and low-latency transmissions by deploying the UAV in close proximity to serve ground users [15]–[18]. Different from a ground base station which suffers from highly scattered Rayleigh fading channels, the UAV can exploit a strong line-of-sight (LoS) channel when it is above a certain altitude and the propagation conditions between the UAV and ground users can be approximated as free space. Furthermore, the UAV can be optimally deployed between the ground users and cloud server to further reduce the transmission and backhaul latency, which is perfectly suitable for latency-sensitive applications. Several papers have addressed the performance of UAV-enabled MEC systems with computing resource constraints [19]–[21]. In [19], a security maximization UAV-enabled MEC framework was proposed by jointly optimizing the UAV location, users' transmit power, UAV jamming power, offloading ratio, UAV computing capacity, and offloading user association. The authors in [20] developed a low-complexity power minimization algorithm by jointly optimizing user association, power control, computation capacity allocation and location planning in a MEC network with multiple UAVs. In [21], the UAV trajectory, user association and user offloading ratio were jointly optimized to minimize the maximum latency in UAV-MEC networks. The performance of a UAV-enabled caching system was investigated in [22]–[24]. In [22], a secure transmission scheme was proposed for a UAV-enabled caching system based on interference alignment. In [23], the UAV location, content caching decision and user association were jointly optimized to maximize the users' quality-of-experience. In [24], the file caching policy, UAV trajectory and file transmission scheduling were jointly optimized in a UAV-enabled network with proactive caching. Notably, no prior works have jointly considered the communication, computing, and caching (3C) performance of UAV systems which is critical for successful low-latency VR delivery, thus motivating this work.

In this paper, we present a novel framework with the aim of minimizing the maximum latency of a UAV-enabled communication, computing and caching VR delivery system as shown in Fig. 1, which consists of one cloud server, one UAV aerial base station (U-BS) equipped with both caching and computing capabilities, and multiple ground VR users with local computing resources. To reduce the traffic burden on the backhaul link and backhaul communication latency, the U-BS caches the most popular input data in its cache container and the data which has not been cached at the U-BS needs to be transmitted from the cloud server via a wireless backhaul link. Moreover, to further reduce latency, the U-BS may choose to process the input data with its computing resource and transmit the projected output data to VR users for display, or send the input data to VR users directly for local computing. We note that compared to [5]–[7] where a VR delivery system was proposed for ground communications, our work exploits the advantages of UAV communications where the latency consumption can be further reduced by optimizing the UAV location, UAV computing capacity allocation and UAV caching policy. In addition, our work which jointly considers the computing and caching capabilities of UAV-enabled systems

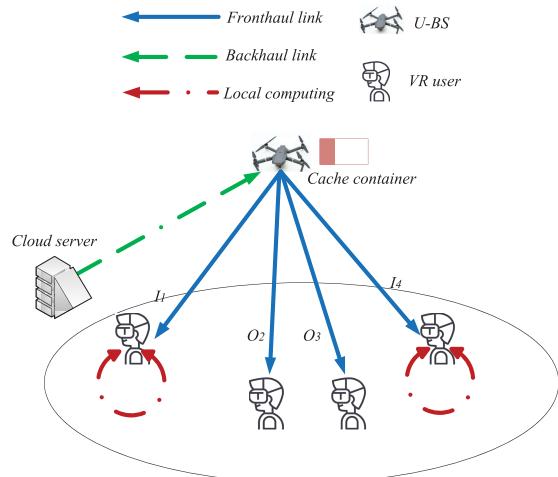


Fig. 1. UAV-enabled communication, computing and caching VR delivery system.

is different from other research on UAV communications such as [19]–[21] and [22]–[24] which solely focused on either computing or caching capabilities, respectively.

The main contributions of this paper are summarized as follows.

- We formulate a maximum latency minimization problem of a UAV-enabled VR delivery system by jointly optimizing the U-BS location, fronthaul and backhaul bandwidth allocation, computing capacity allocation, data caching policy and computing offloading policy subject to computing, caching and power constraints.
- To solve the non-convex optimization problem, we first apply the block coordinate descent (BCD) method to decouple the original optimization problem into six subproblems and propose a low-complexity algorithm to solve each subproblem alternately. We solve the U-BS location subproblem by applying a successive convex approximation (SCA) on the U-BS data rate. Then, we apply Lagrangian dual decomposition method to efficiently solve the bandwidth and computing capacity allocation subproblems. Finally, we obtain efficient closed-form solutions for the caching and computing policy subproblems.
- Simulation results show that our proposed algorithm achieves a lower latency compared to benchmark strategies and highlight a tradeoff between latency and the primary resource requirements of communication, computing and caching.

The rest of this paper is organized as follows. Section II introduces the UAV-enabled communication, computing and caching VR delivery system model and formulates the joint optimization problem. In Section III, we propose an efficient iterative algorithm to minimize the maximum latency consumption. The effectiveness of our proposed solution is shown through simulation results in Section IV. Finally, we conclude the paper in Section V.

## II. SYSTEM MODEL

Fig. 1 depicts our proposed UAV-enabled communication, computing and caching VR delivery system with  $N$  ground VR users, one U-BS and one ground cloud server, where the

TABLE I  
TABLE OF NOTATIONS

Notation	Description
$\mathcal{N}$	Set of VR users
$\theta_i$	Fraction of fronthaul bandwidth allocated to the $i$ -th VR user
$\eta_i$	Fraction of backhaul bandwidth allocated for transmitting $I_i$
$f_i$	Computing capacity of U-BS assigned to the $i$ -th VR user
$c_i$	Caching policy variable
$a_i$	Computing policy variable
$p_u, p_c$	Transmit power for the U-BS and cloud server
$\sigma^2$	Noise spectral density
$\beta_0$	Reference channel power gain
$\mathbf{y}, \mathbf{w}_i, \mathbf{v}$	Horizontal location of the U-BS, the $i$ -th VR user and cloud server
$H_u$	Altitude of U-BS
$r_i$	Fronthaul data rate at the $i$ -th VR user
$r_i^b$	Backhaul data rate for transmitting $I_i$
$I_i, O_i$	Input and output data size of the $i$ -th VR user
$\alpha$	Ratio of size between $O_i$ and $I_i$
$F_i$	CPU cycles for computing data $I_i$
$f_i^{local}$	Local computing capacity at the $i$ -th VR user

set of VR users is denoted as  $\mathcal{N} = \{1, 2, \dots, N\}$ . We consider that the U-BS has caching and computing capabilities which enable it to cache the data requested by each VR user from the cloud server via wireless backhaul and compute the data, respectively. Each VR user with computing capability is able to compute the data locally. We assume that all devices are equipped with a single antenna for transmitting or receiving. Due to the long distance and blockages from the cloud server to the VR users, the direct links between them are not applicable.

#### A. Computing Model

We assume that the  $i$ -th VR user has the computationally intensive task  $U_i$  to be executed as follows [13]

$$U_i = (I_i, O_i, F_i), \quad \forall i \in \mathcal{N}, \quad (1)$$

where  $I_i$  is the input data in bits of the VR video required by the  $i$ -th user which is available in the remote cloud server, and may or may not be cached at the U-BS,  $O_i = \alpha I_i$  is the output data in bits after being processed at the U-BS or locally with  $\alpha \geq 2$  as the ratio of size between  $O_i$  and  $I_i$  [13], and  $F_i$  is the number of CPU cycles for computing one bit of input data  $I_i$ .

We consider VR projection and rendering in our system and define  $a_i = \{0, 1\}, \forall i \in \mathcal{N}$ , as the computing policy variable where  $a_i = 1$  indicates that the input data  $I_i$  required by the  $i$ -th VR user will be projected at the U-BS. Thus, the U-BS processes the input data and transmits the output data  $O_i$  to VR users for display. On the other hand,  $a_i = 0$  indicates that the  $i$ -th VR user decides to compute its data locally, but this user also needs to receive the input data,  $I_i$ , from the U-BS for calculation. Thus, the fronthaul transmission latency between the U-BS and each VR user is jointly decided by the computing policy, data size, and transmission rate, which is given by

$$t_i^{tr} = a_i \cdot \frac{O_i}{r_i} + (1 - a_i) \cdot \frac{I_i}{r_i}, \quad \forall i \in \mathcal{N}, \quad (2)$$

where  $r_i$  is the transmission rate between the U-BS and the  $i$ -th VR user which is shown in (10). The first term in the right-hand-side (RHS) of (2) shows that if the data is computed at the U-BS, the output data  $O_i$  after being processed will be transmitted from U-BS to the VR user and the second term means that if the data is computed locally, the U-BS transmits the input data  $I_i$  to the VR user for calculation.

The computing latency which depends on computing policy, data size, computing capacity, and the required CPU cycles of the computing data, is given by

$$t_i^c = a_i \cdot \frac{I_i \cdot F_i}{f_i} + (1 - a_i) \cdot \frac{I_i \cdot F_i}{f_i^{local}}, \quad \forall i \in \mathcal{N}, \quad (3)$$

where  $f_i^{local}$  is the local computing capacity at the  $i$ -th VR user and  $f_i$  is the computing capacity of the U-BS assigned to compute the data requested by the  $i$ -th VR user, which is constrained by a maximum computing capacity given by

$$\sum_{i=1}^N a_i f_i \leq f_{max}. \quad (4)$$

We note that if the  $i$ -th VR user decides to locally compute its data and  $a_i = 0$ , the U-BS will not allocate any computing capacity to this VR user and  $f_i = 0$ . We set the first term in the RHS of (3) to zero when  $a_i = 0$  and  $f_i = 0$ .

We model the power consumption at the U-BS for computing the input data requested by the  $i$ -th VR user as [19]

$$p_i^c = \kappa f_i^3, \quad \forall i \in \mathcal{N}, \quad (5)$$

where  $\kappa$  is the effective switched capacitance on the chip. The total power consumption at the U-BS which consists of transmit power,  $p_u$ , and computing power should be limited by a maximum budget  $p_{max}$ , which is given by

$$p_u + \sum_{i=1}^N a_i \kappa f_i^3 \leq p_{max}. \quad (6)$$

### B. Caching Model

We define  $c_i = \{0, 1\}, \forall i \in \mathcal{N}$  as the caching policy variable where  $c_i = 1$  represents that the input data requested by the  $i$ -th VR user has been cached in the U-BS and  $c_i = 0$  otherwise. We note that if the U-BS has cached the input data requested by the  $i$ -th VR user, it can apply the data directly from its cache container. Otherwise, the input data has to be transmitted from the cloud server to the U-BS and the corresponding backhaul latency is given by

$$t_i^b = (1 - c_i) \cdot \frac{I_i}{r_i^b}, \quad \forall i \in \mathcal{N}, \quad (7)$$

where  $r_i^b$  representing the backhaul rate for transmitting  $I_i$  is given in (12).

Since different portions of the VR video are viewed by different VR users based on their geographical locations, we assume that the input data required by each VR user is different from each other. Note that the caching storage at the U-BS should be bounded by a maximum budget  $c_{max}$ , which is given by

$$\sum_{i=1}^N c_i I_i \leq c_{max}. \quad (8)$$

### C. Communication Model

Assume that the coordinate of the  $i$ -th VR user is denoted by  $\mathbf{w}_i = (x_i, y_i)^T \in \mathbb{R}^{2 \times 1}, \forall i \in \mathcal{N}$ . The U-BS is fixed at altitude  $H_u$ , which is the minimum altitude required by regulations to avoid building obstacles, and its horizontal location is denoted by  $\mathbf{y} = (x_u, y_u)^T \in \mathbb{R}^{2 \times 1}$ . For the air-to-ground channel, we adopt a simple channel model where the channel power gains are dominated by the LoS links. Then, the channel power gain between the U-BS and the  $i$ -th VR user is given as [25], [26]

$$h_i = \frac{\beta_0}{\|\mathbf{y} - \mathbf{w}_i\|^2 + H_u^2}, \quad \forall i \in \mathcal{N}, \quad (9)$$

where  $\beta_0$  denotes the channel power gain at the reference distance of one meter.

Define  $\theta_i \geq 0, \forall i \in \mathcal{N}$  as the fronthaul bandwidth allocation factor which represents the fraction of fronthaul bandwidth that the U-BS allocates to the  $i$ -th VR user. The achievable data rate at the  $i$ -th VR user is denoted by  $r_i$  in bits/second (bps), which is expressed as

$$r_i = \theta_i B \log_2 \left( 1 + \frac{p_u h_i}{\theta_i B \sigma^2} \right), \quad \forall i \in \mathcal{N}, \quad (10)$$

where  $p_u$  is the transmit power at the U-BS,  $B$  is the total fronthaul bandwidth, and  $\sigma^2$  is the noise spectral density.

Assume that the coordinate of the cloud server is denoted by  $\mathbf{v} = (x_c, y_c)^T \in \mathbb{R}^{2 \times 1}$ . Then, the channel power gain between the cloud server and the U-BS is given as [25], [26]

$$h_b = \frac{\beta_0}{\|\mathbf{y} - \mathbf{v}\|^2 + H_u^2}. \quad (11)$$

We define  $\mathcal{N}_{uncached} = \{i | c_i = 0, \forall i \in \mathcal{N}\}$  as the set of VR users whose input data has not been cached in the U-BS and  $\eta_i \geq 0, \forall i \in \mathcal{N}_{uncached}$  as the backhaul

bandwidth allocation factor which represents the fraction of backhaul bandwidth that the cloud allocates to transmit the input data  $I_i$  which has not been cached in the U-BS. We note that for the  $i$ -th VR user whose requested input data  $I_i$  has been cached in the U-BS, i.e.,  $c_i = 1$ , the cloud will not allocate any backhaul bandwidth for transmitting  $I_i$  and  $\eta_i = 0, \forall i \in \mathcal{N} / \mathcal{N}_{uncached}$ . The achievable backhaul data rate for transmitting  $I_i$  is denoted by  $r_i^b$  in bits/second (bps), which is expressed as

$$r_i^b = \eta_i B_{back} \log_2 \left( 1 + \frac{p_c h_b}{\eta_i B_{back} \sigma^2} \right), \quad \forall i \in \mathcal{N}_{uncached}, \quad (12)$$

where  $p_c$  is the transmit power at the cloud server and  $B_{back}$  is the total backhaul bandwidth.

According to (2), (3), (7), (10), and (12), the total latency to complete the task at each VR user is given by

$$\begin{aligned} t_i &= t_i^{tr} + t_i^c + t_i^b \\ &= a_i \left( \frac{O_i}{r_i} + \frac{I_i F_i}{f_i} \right) + (1 - a_i) \left( \frac{I_i F_i}{f_i^{local}} + \frac{I_i}{r_i} \right) \\ &\quad + \frac{(1 - c_i) I_i}{r_i^b}, \quad \forall i \in \mathcal{N}. \end{aligned} \quad (13)$$

### D. Problem Formulation

We note that the satisfaction of VR experience among all users is dominated by the user who experiences the worst latency. To achieve the fairness among all VR users, we formulate an optimization problem aimed at minimizing the maximum latency among all VR users subject to computing capacity, caching storage and total power constraints. We jointly optimize the U-BS location  $\mathbf{y} = \{(x_u, y_u)^T\}$ , fronthaul bandwidth allocation  $\boldsymbol{\theta} = \{\theta_i, \forall i \in \mathcal{N}\}$ , backhaul bandwidth allocation  $\boldsymbol{\eta} = \{\eta_i, \forall i \in \mathcal{N}_{uncached}\}$ , computing capacity allocation  $\mathbf{f} = \{f_i, \forall i \in \mathcal{N}\}$ , data caching policy  $\mathbf{c} = \{c_i, \forall i \in \mathcal{N}\}$ , and computing policy  $\mathbf{a} = \{a_i, \forall i \in \mathcal{N}\}$ . The optimization problem can be formulated as

$$\min_{\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\eta}, \mathbf{f}, \mathbf{c}, \mathbf{a}} \max_{i \in \mathcal{N}} t_i \quad (14a)$$

$$\text{s.t.} \quad \sum_{i=1}^N a_i f_i \leq f_{max} \quad (14b)$$

$$\sum_{i=1}^N c_i I_i \leq c_{max} \quad (14c)$$

$$p_u + \sum_{i=1}^N a_i \kappa f_i^3 \leq p_{max} \quad (14d)$$

$$a_i = \{0, 1\}, \quad \forall i \in \mathcal{N} \quad (14e)$$

$$c_i = \{0, 1\}, \quad \forall i \in \mathcal{N} \quad (14f)$$

$$\sum_{i=1}^N \theta_i \leq 1 \quad (14g)$$

$$\sum_{i=1}^N \eta_i \leq 1. \quad (14h)$$

Define an auxiliary variable  $T \triangleq \max_{i \in \mathcal{N}} t_i$  as the maximum latency among all VR users, we can reformulate the original

optimization problem as

$$\min_{\mathbf{y}, \boldsymbol{\theta}, \boldsymbol{\eta}, \mathbf{f}, \mathbf{c}, T} T \quad (15a)$$

$$\text{s.t. } a_i \left( \frac{O_i}{r_i} + \frac{I_i F_i}{f_i} \right) + (1 - a_i) \left( \frac{I_i F_i}{f_i^{local}} + \frac{I_i}{r_i} \right) + \frac{(1 - c_i) I_i}{r_i^b} \leq T, \quad \forall i \in \mathcal{N} \quad (15b)$$

(14b) – (14h),

where the newly defined constraint (15b) is based on the intrinsic limitation that the latency consumption of each user should be less than its maximum value. Although reformulated, Problem (15) is still a non-convex optimization problem due to the following reasons. First, the optimizing variables for computing policy  $\mathbf{a}$  and data caching policy  $\mathbf{c}$  are binary integers. Second, even with given  $\mathbf{a}$  and  $\mathbf{c}$ , (15b) is still a non-convex constraint with respect to U-BS location  $\mathbf{y}$ . Therefore, the main challenge that we will address in the following section is to develop an efficient algorithm to solve the latency optimization problem in (15).

### III. PROPOSED LATENCY MINIMIZATION ALGORITHM

In this section, we detail our proposed latency minimization algorithm for UAV-enabled VR delivery systems. To solve Problem (15), we apply the BCD method which alternately optimizes one block of optimization variable in each iteration while keeping other blocks of optimization variables fixed to obtain a high-quality suboptimal solution [15]. Therefore, we can decouple the original optimization problem into six subproblems to solve the U-BS location  $\mathbf{y}$ , fronthaul bandwidth allocation  $\boldsymbol{\theta}$ , backhaul bandwidth allocation  $\boldsymbol{\eta}$ , computing capacity allocation  $\mathbf{f}$ , data caching policy  $\mathbf{c}$ , and computing policy  $\mathbf{a}$  in an iterative manner.

#### A. U-BS Location Subproblem

For any given  $\boldsymbol{\theta}, \boldsymbol{\eta}, \mathbf{f}, \mathbf{c}$ , and  $\mathbf{a}$ , the U-BS location of Problem (15) can be optimized by solving the following problem

$$\min_{\mathbf{y}, T} T \quad (16a)$$

$$\text{s.t. } \frac{a_i O_i + (1 - a_i) I_i}{\theta_i B \log_2 \left( 1 + \frac{\zeta_i}{\|\mathbf{y} - \mathbf{w}_i\|^2 + H_u^2} \right)} + \frac{(1 - c_i) I_i}{\eta_i B_{back} \log_2 \left( 1 + \frac{\gamma_i}{\|\mathbf{y} - \mathbf{v}\|^2 + H_u^2} \right)} \leq T - \rho_i, \quad \forall i \in \mathcal{N}, \quad (16b)$$

where the constraint (16b) corresponds to (15b), and all the other constraints in (15) are not applicable. In (16), we define  $\zeta_i = \frac{p_u \beta_0}{\theta_i B \sigma^2}$ ,  $\gamma_i = \frac{p_c \beta_0}{\eta_i B_{back} \sigma^2}$ , and  $\rho_i = a_i \frac{I_i F_i}{f_i} + (1 - a_i) \frac{I_i}{f_i^{local}}$ . Note that (16) is a non-convex optimization problem and the non-convexity arises from the logarithm terms. In the following, we first introduce slack variables  $\boldsymbol{\epsilon} \triangleq \{\epsilon_i, \forall i \in \mathcal{N}\}$  and  $\boldsymbol{\omega} \triangleq \{\omega_i, \forall i \in \mathcal{N}_{uncached}\}$ , and reformulate the U-BS location subproblem as

$$\min_{\mathbf{y}, \boldsymbol{\epsilon}, \boldsymbol{\omega}, T} T \quad (17a)$$

$$\text{s.t. } \frac{a_i O_i + (1 - a_i) I_i}{\epsilon_i} + \frac{(1 - c_i) I_i}{\omega_i} \leq T - \rho_i, \quad \forall i \in \mathcal{N} \quad (17b)$$

$$\underbrace{\theta_i B \log_2 \left( 1 + \frac{\zeta_i}{\|\mathbf{y} - \mathbf{w}_i\|^2 + H_u^2} \right)}_{\mathcal{I}_i} \geq \epsilon_i, \quad \forall i \in \mathcal{N} \quad (17c)$$

$$\underbrace{\eta_i B_{back} \log_2 \left( 1 + \frac{\gamma_i}{\|\mathbf{y} - \mathbf{v}\|^2 + H_u^2} \right)}_{\mathcal{Z}_i} \geq \omega_i, \quad \forall i \in \mathcal{N}_{uncached}. \quad (17d)$$

We note that the constraint (17b) is convex now and the non-convexity of Problem (17) arises from constraints (17c) and (17d). Next, we adopt the SCA technique, where the original function can be approximated by a more tractable expression at a given local point in each iteration [15], [19]. We note that  $\mathcal{I}_i$  is convex with respect to  $\|\mathbf{y} - \mathbf{w}_i\|^2$ , thus, a concave lower bound expression  $\mathcal{I}_i^{lb}$  can be derived by applying the first-order Taylor expansion with given U-BS location  $\mathbf{y}[m]$  in the  $m$ -th iteration, which is given by

$$\begin{aligned} \mathcal{I}_i^{lb} = & \theta_i B \log_2 \left( 1 + \frac{\zeta_i}{\|\mathbf{y}[m] - \mathbf{w}_i\|^2 + H_u^2} \right) \\ & - \frac{\theta_i B \zeta_i (\|\mathbf{y} - \mathbf{w}_i\|^2 - \|\mathbf{y}[m] - \mathbf{w}_i\|^2)}{(\|\mathbf{y}[m] - \mathbf{w}_i\|^2 + H_u^2 + \zeta_i)(\|\mathbf{y}[m] - \mathbf{w}_i\|^2 + H_u^2) \ln 2}. \end{aligned} \quad (18)$$

We apply a similar approach on  $\mathcal{Z}_i$  and the corresponding concave lower bound  $\mathcal{Z}_i^{lb}$  is given by

$$\begin{aligned} \mathcal{Z}_i^{lb} = & \eta_i B_{back} \log_2 \left( 1 + \frac{\gamma_i}{\|\mathbf{y}[m] - \mathbf{v}\|^2 + H_u^2} \right) \\ & - \frac{\eta_i B_{back} \gamma_i (\|\mathbf{y} - \mathbf{v}\|^2 - \|\mathbf{y}[m] - \mathbf{v}\|^2)}{(\|\mathbf{y}[m] - \mathbf{v}\|^2 + H_u^2 + \gamma_i)(\|\mathbf{y}[m] - \mathbf{v}\|^2 + H_u^2) \ln 2}. \end{aligned} \quad (19)$$

With given U-BS location  $\mathbf{y}[m]$  and the lower bound expressions in (18) and (19), the U-BS location subproblem can be solved as

$$\min_{\mathbf{y}, \boldsymbol{\epsilon}, \boldsymbol{\omega}, T} T \quad (20a)$$

$$\text{s.t. } \frac{a_i O_i + (1 - a_i) I_i}{\epsilon_i} + \frac{(1 - c_i) I_i}{\omega_i} \leq T - \rho_i, \quad \forall i \in \mathcal{N} \quad (20b)$$

$$\mathcal{I}_i^{lb} \geq \epsilon_i, \quad \forall i \in \mathcal{N} \quad (20c)$$

$$\mathcal{Z}_i^{lb} \geq \omega_i, \quad \forall i \in \mathcal{N}_{uncached}. \quad (20d)$$

We note that Problem (20) is a convex optimization problem and it can be efficiently solved by utilizing mathematical optimization software with the polynomial complexity [27].

#### B. Fronthaul Bandwidth Allocation Subproblem

For any given  $\mathbf{y}, \boldsymbol{\eta}, \mathbf{f}, \mathbf{c}$ , and  $\mathbf{a}$ , the fronthaul bandwidth allocation of Problem (15) can be optimized by solving the following problem

$$\min_{\boldsymbol{\theta}, T} T \quad (21a)$$

$$\text{s.t. } \frac{\chi_i}{T - \nu_i} \leq \theta_i B \log_2 \left( 1 + \frac{p_u h_i}{\theta_i B \sigma^2} \right), \quad \forall i \in \mathcal{N} \quad (21b)$$

$$\sum_{i=1}^N \theta_i \leq 1 \quad (21c)$$

$$\theta_i \geq 0, \quad \forall i \in \mathcal{N} \quad (21d)$$

$$T \geq \nu_i, \quad \forall i \in \mathcal{N}, \quad (21e)$$

where  $\nu_i = a_i \frac{I_i F_i}{f_i} + (1 - a_i) \frac{I_i F_i}{f_{local}} + \frac{(1 - c_i) I_i}{r_i^b}$  and  $\chi_i = a_i O_i + (1 - a_i) I_i$ . We define  $\theta_i B \log_2 \left( 1 + \frac{p_u h_i}{\theta_i B \sigma^2} \right) \triangleq 0$  when  $\theta_i = 0, \forall i \in \mathcal{N}$ , such that the RHS of (21b) is continuous with respect to  $\theta_i$  over the whole domain. We analyze the convexity of Problem (21) in the following lemma.

*Lemma 1:* Problem (21) is a convex problem.

*Proof:* It can be easily noted that (21a), (21c), (21d), and (21e) are convex terms due to their linearity. Therefore, proving Lemma 1 is equivalent to proving that the constraint (21b) is convex. To show this, we define  $g(x) = x \log_2 \left( 1 + \frac{1}{x} \right), x > 0$ , and we have

$$\frac{\partial^2 g}{\partial x^2} = -\frac{1}{x(x+1)^2 \ln 2} < 0, \quad \forall x > 0, \quad (22)$$

which indicates that  $g(x)$  is a concave function. Thus, we can conclude that the RHS of (21b) is a concave term with respect to  $\theta_i$ . Moreover, We note that the left-hand-side (LHS) of (21b) is a convex term with respect to  $T$ . Therefore, we show that the constraint (21b) is convex and prove that Problem (21) is a convex problem.  $\square$

Next, we apply the Lagrangian dual decomposition method to solve this convex problem. It can be verified that the Slater's condition is satisfied for Problem (21), which indicates that the duality gap between (21) and its dual problem is zero [27]. The partial Lagrangian function of Problem (21) is given by

$$\begin{aligned} \mathcal{L}(T, \boldsymbol{\theta}, \boldsymbol{\mu}, \iota) = & T + \sum_{i=1}^N \frac{\mu_i \chi_i}{T - \nu_i} \\ & + \sum_{i=1}^N \left[ \iota \theta_i - \mu_i \theta_i B \log_2 \left( 1 + \frac{p_u h_i}{\theta_i B \sigma^2} \right) \right] - \iota, \end{aligned} \quad (23)$$

where  $\boldsymbol{\mu} = \{\mu_i, \forall i \in \mathcal{N}\}$  and  $\iota$  are Lagrangian multipliers associated with constraints (21b) and (21c), respectively. The boundary constraints (21d) and (21e) will be absorbed into the optimal solution in the following. The dual function is given by

$$f(\boldsymbol{\mu}, \iota) = \min_{T, \boldsymbol{\theta}} \mathcal{L}(T, \boldsymbol{\theta}, \boldsymbol{\mu}, \iota) \quad (24a)$$

$$s.t. \quad \theta_i \geq 0, T \geq \nu_i, \quad \forall i \in \mathcal{N}, \quad (24b)$$

and the dual problem of (21) is given by

$$\max_{\boldsymbol{\mu}, \iota} f(\boldsymbol{\mu}, \iota) \quad (25a)$$

$$s.t. \quad \boldsymbol{\mu} \succeq 0, \iota \geq 0. \quad (25b)$$

To derive the primal optimal solution of Problem (21), we apply the Lagrange duality and derive  $f(\boldsymbol{\mu}, \iota)$  by solving Problem (24). We note that with given dual variables  $\boldsymbol{\mu}$  and  $\iota$ , Problem (24) can be decomposed into  $N + 1$  independent subproblems where one subproblem is for optimizing  $T$  and

the other  $N$  subproblems are for optimizing  $\theta_i, \forall i \in \mathcal{N}$ . The subproblem for optimizing  $T$  can be formulated as

$$\min_T T + \sum_{i=1}^N \frac{\mu_i \chi_i}{T - \nu_i} \quad (26a)$$

$$s.t. \quad T \geq \nu_i, \quad \forall i \in \mathcal{N}. \quad (26b)$$

By setting the first-order derivative of (26a) with respect to  $T$  to zero, we observe that the optimal  $T$  should satisfy

$$T = \left\{ T \mid \sum_{i=1}^N \frac{\mu_i \chi_i}{(T - \nu_i)^2} = 1, T \geq \nu_i \right\}, \quad (27)$$

which can be found by applying the bisection method.

Moreover, the subproblem for optimizing  $\theta_i, \forall i \in \mathcal{N}$  can be formulated as

$$\min_{\theta_i} \iota \theta_i - \mu_i \theta_i B \log_2 \left( 1 + \frac{p_u h_i}{\theta_i B \sigma^2} \right) \quad (28a)$$

$$s.t. \quad \theta_i \geq 0, \quad \forall i \in \mathcal{N}. \quad (28b)$$

By setting the first-order derivative of (28a) with respect to  $\theta_i$  to zero, we obtain the closed-form expression of the optimal bandwidth allocation as

$$\theta_i = \left[ \frac{p_u h_i}{B \sigma^2} \left( -\frac{1}{\mathcal{W} \left( -\frac{1}{\exp(1 + \frac{\iota \ln 2}{\mu_i B})} \right)} - 1 \right)^{-1} \right]^+, \quad (29)$$

where  $[x]^+ = \max\{x, 0\}$  and  $\mathcal{W}(x)$  is the Lambert function, which is defined as the inverse function of  $f(x) = x \exp(x)$ .

The value of dual variables  $\boldsymbol{\mu}$  and  $\iota$  can be determined by the sub-gradient method. The updating procedure can be given by

$$\mu_i = \left[ \mu_i + \phi \left( \frac{\chi_i}{T - \nu_i} - \theta_i B \log_2 \left( 1 + \frac{p_u h_i}{\theta_i B \sigma^2} \right) \right) \right]^+, \quad \forall i \in \mathcal{N} \quad (30a)$$

$$\iota = \left[ \iota + \phi \left( \sum_{i=1}^N \theta_i - 1 \right) \right]^+, \quad (30b)$$

where  $\phi > 0$  is a dynamic step-size sequence, which can be selected by using the typical self-adaptive scheme [18].

We note that in the primal problem of (21), the optimal  $T$  and  $\boldsymbol{\theta}$  can be derived by solving (27) and (29), respectively. Moreover, in the dual problem of (21), the optimal dual variables  $\boldsymbol{\mu}$  and  $\iota$  can be found by solving (30a) and (30b), respectively. The details for obtaining the optimal solution to Problem (21) are summarized in Algorithm 1. We note that Problem (24) has been decomposed into  $N + 1$  subproblems. To solve the subproblem for optimizing  $T$ , the complexity of solving (27) via the bisection method is  $\mathcal{O}(\log_2(1/\epsilon))$  with  $\epsilon$  being the iterative accuracy. To solve each of the  $N$  subproblems, since the closed-form solution has been derived in (29), the complexity for these  $N$  subproblems is  $\mathcal{O}(N)$ . Moreover, the complexity of updating dual variables is  $\mathcal{O}(N)$  according to (30a) and (30b). As a result, the total complexity of Algorithm 1 is  $\mathcal{O}(L_1 L_2 N^2 \log_2(1/\epsilon))$ , where  $L_1$  is the number of iterations for outer layer in Algorithm 1 and  $L_2$  is the number of iterations via the dual method of solving Problem (21).

**Algorithm 1** Fronthaul Bandwidth Allocation Algorithm for Solving Problem (21)

- 
- 1: Initialize  $\mu$  and  $\iota$ .
  - 2: **repeat**
  - 3:   Obtain the optimal  $T$  and  $\theta$  by solving (27) and (29), respectively;
  - 4:   Update the Lagrangian multipliers  $\mu$  and  $\iota$  by solving (30a) and (30b), respectively;
  - 5: **until** The objective function in (21a) converges.
- 

**C. Backhaul Bandwidth Allocation Subproblem**

For any given  $\mathbf{y}, \theta, \mathbf{f}, \mathbf{c}$ , and  $\mathbf{a}$ , the backhaul bandwidth allocation of Problem (15) can be optimized by solving the following problem

$$\min_{\eta, T} T \quad (31a)$$

$$\text{s.t. } \frac{I_i}{T - z_i} \leq \eta_i B_{back} \log_2 \left( 1 + \frac{p_c h_b}{\eta_i B_{back} \sigma^2} \right), \quad \forall i \in \mathcal{N}_{uncached} \quad (31b)$$

$$\sum_{i=1}^{N_{uncached}} \eta_i \leq 1 \quad (31c)$$

$$\eta_i \geq 0, \quad \forall i \in \mathcal{N}_{uncached} \quad (31d)$$

$$T \geq z_i, \quad \forall i \in \mathcal{N}_{uncached}, \quad (31e)$$

where  $z_i = a_i \left( \frac{O_i}{r_i} + \frac{I_i F_i}{f_i} \right) + (1 - a_i) \left( \frac{I_i F_i}{f_i^{local}} + \frac{I_i}{r_i} \right)$ . We set  $\eta_i B_{back} \log_2 \left( 1 + \frac{p_c h_b}{\eta_i B_{back} \sigma^2} \right) \triangleq 0$  when  $\eta_i = 0, \forall i \in \mathcal{N}_{uncached}$ , such that the RHS of (31b) is continuous with respect to  $\eta_i$  over the whole domain. We note that Problem (31) is a convex problem and the proof is similar to that of (21) in Subsection III-B.

As such, we can apply the Lagrangian dual decomposition method to solve Problem (31). We denote  $\lambda = \{\lambda_i, \forall i \in \mathcal{N}_{uncached}\}$  and  $\varsigma$  as Lagrangian multipliers associated with constraints (31b) and (31c), respectively. The boundary constraints (31d) and (31e) will be absorbed into the optimal solution in the following. The partial Lagrangian function of Problem (31) is given by

$$\begin{aligned} \mathcal{L}(T, \eta, \lambda, \varsigma) = & T + \sum_{i=1}^{N_{uncached}} \frac{\lambda_i I_i}{T - z_i} \\ & + \sum_{i=1}^{N_{uncached}} \left[ \varsigma \eta_i - \lambda_i \eta_i B_{back} \log_2 \left( 1 + \frac{p_c h_b}{\eta_i B_{back} \sigma^2} \right) \right] - \varsigma. \end{aligned} \quad (32)$$

The dual function is given by

$$f(\lambda, \varsigma) = \min_{T, \eta} \mathcal{L}(T, \eta, \lambda, \varsigma) \quad (33a)$$

$$\text{s.t. } \eta_i \geq 0, T \geq z_i, \quad \forall i \in \mathcal{N}_{uncached}, \quad (33b)$$

and the dual problem of (31) is given by

$$\max_{\lambda, \varsigma} f(\lambda, \varsigma) \quad (34a)$$

$$\text{s.t. } \lambda \succeq 0, \varsigma \geq 0. \quad (34b)$$

To derive the primal optimal solution of Problem (31), we apply the Lagrange duality method and derive  $f(\lambda, \varsigma)$  by solving Problem (33). We note that with given dual variables  $\lambda$  and  $\varsigma$ , Problem (33) can be decomposed into  $N_{uncached} + 1$  independent subproblems where one subproblem is for optimizing  $T$  and the other  $N_{uncached}$  subproblems are for optimizing  $\eta_i, \forall i \in \mathcal{N}_{uncached}$ . The subproblem for optimizing  $T$  can be formulated as

$$\min_T T + \sum_{i=1}^{N_{uncached}} \frac{\lambda_i I_i}{T - z_i} \quad (35a)$$

$$\text{s.t. } T \geq z_i, \quad \forall i \in \mathcal{N}_{uncached}. \quad (35b)$$

By setting the first-order derivative of (35a) with respect to  $T$  to zero, we find that the optimal  $T$  should satisfy

$$T = \left\{ T \mid \sum_{i=1}^{N_{uncached}} \frac{\lambda_i I_i}{(T - z_i)^2} = 1, T \geq z_i \right\}, \quad (36)$$

which can be solved by applying the bisection search method.

Moreover, the subproblem for optimizing  $\eta_i, \forall i \in \mathcal{N}_{uncached}$  can be formulated as

$$\min_{\eta_i} \varsigma \eta_i - \lambda_i \eta_i B_{back} \log_2 \left( 1 + \frac{p_c h_b}{\eta_i B_{back} \sigma^2} \right) \quad (37a)$$

$$\text{s.t. } \eta_i \geq 0, \quad \forall i \in \mathcal{N}_{uncached}. \quad (37b)$$

By setting the first-order derivative of (37a) with respect to  $\eta_i$  to zero, we obtain the closed-form expression of the optimal backhaul bandwidth allocation as

$$\eta_i = \left[ \frac{p_c h_b}{B_{back} \sigma^2} \left( -\frac{1}{\mathcal{W} \left( -\frac{1}{\exp(1 + \frac{\varsigma \ln 2}{\lambda_i B_{back}})} \right)} - 1 \right)^{-1} \right]^+. \quad (38)$$

The value of dual variables  $\lambda$  and  $\varsigma$  can be determined by the sub-gradient method. The updating procedure is given by

$$\lambda_i = \left[ \lambda_i + \phi \left( \frac{I_i}{T - z_i} - \eta_i B_{back} \log_2 \left( 1 + \frac{p_c h_b}{\eta_i B_{back} \sigma^2} \right) \right) \right]^+ \quad \forall i \in \mathcal{N}_{uncached} \quad (39a)$$

$$\varsigma = \left[ \varsigma + \phi \left( \sum_{i=1}^{N_{uncached}} \eta_i - 1 \right) \right]^+. \quad (39b)$$

The procedures for obtaining the optimal solution to Problem (31) is summarized in Algorithm 2. Similar to the complexity analysis in Subsection III-B, we note that the total complexity of Algorithm 2 is  $\mathcal{O}(L_3 L_4 N^2 \log_2(1/\epsilon))$ , where  $L_3$  is the number of iterations for outer layer in Algorithm 2 and  $L_4$  is the number of iterations via the dual method of solving Problem (31).

**D. Computing Capacity Allocation Subproblem**

We define  $\mathcal{N}_{as} = \{i | a_i = 1, \forall i \in \mathcal{N}\}$  as the set of VR users who choose to compute their input data at the U-BS. Note that for the VR users who choose to self-execute their tasks, the U-BS will not allocate computing capacity to them and  $f_i = 0, \forall i \in \mathcal{N}/\mathcal{N}_{as}$ . For any given  $\mathbf{y}, \theta, \eta, \mathbf{c}$ , and

**Algorithm 2** Backhaul Bandwidth Allocation Algorithm for Solving Problem (31)

- 
- 1: Initialize  $\lambda$  and  $\varsigma$ .
  - 2: **repeat**
  - 3:   Obtain the optimal  $T$  and  $\eta$  by solving (36) and (38), respectively;
  - 4:   Update the Lagrangian multipliers  $\lambda$  and  $\varsigma$  by solving (39a) and (39b), respectively;
  - 5: **until** The objective function in (31a) converges.
- 

a, the computing capacity allocation of Problem (15) can be optimized by solving the following problem

$$\min_{\mathbf{f}, T} T \quad (40a)$$

$$\text{s.t. } \frac{I_i F_i}{T - \varpi_i} \leq f_i, \quad \forall i \in \mathcal{N}_{as} \quad (40b)$$

$$\sum_{i=1}^{N_{as}} \kappa f_i^3 \leq p_{max} - p_u \quad (40c)$$

$$\sum_{i=1}^{N_{as}} f_i \leq f_{max} \quad (40d)$$

$$f_i \geq 0, \quad \forall i \in \mathcal{N}_{as} \quad (40e)$$

$$T \geq \varpi_i, \quad \forall i \in \mathcal{N}_{as}, \quad (40f)$$

where  $\varpi_i = a_i \frac{Q_i}{r_i} + (1 - a_i) \left( \frac{I_i F_i}{f_i^{local}} + \frac{I_i}{r_i} \right) + \frac{(1 - c_i) I_i}{r_i^b}$ . We note that Problem (40) is a convex optimization problem since the objective function and all constraints are convex and it can be effectively solved via the Lagrangian dual decomposition method.

We denote  $\tau = \{\tau_i, \forall i \in \mathcal{N}_{as}\}$ ,  $\delta$ , and  $\xi$  as the Lagrangian multipliers associated with constraints (40b), (40c) and (40d), respectively. The boundary constraints (40e) and (40f) will be absorbed into the optimal solution in the following. The partial Lagrangian function of Problem (40) is given by

$$\begin{aligned} \mathcal{L}(T, \mathbf{f}, \tau, \delta, \xi) = & T + \sum_{i=1}^{N_{as}} \frac{\tau_i I_i F_i}{T - \varpi_i} + \sum_{i=1}^{N_{as}} (\delta \kappa f_i^3 + \xi f_i - \tau_i f_i) \\ & + \delta p_u - \delta p_{max} - \xi f_{max}. \end{aligned} \quad (41)$$

The dual function is given by

$$f(\tau, \delta, \xi) = \min_{T, \mathbf{f}} \mathcal{L}(T, \mathbf{f}, \tau, \delta, \xi) \quad (42a)$$

$$\text{s.t. } f_i \geq 0, T \geq \varpi_i, \quad \forall i \in \mathcal{N}_{as}, \quad (42b)$$

and the dual problem of (40) is given by

$$\max_{\tau, \delta, \xi} f(\tau, \delta, \xi) \quad (43a)$$

$$\text{s.t. } \tau \succeq 0, \delta \geq 0, \xi \geq 0. \quad (43b)$$

To derive the primal optimal solution of Problem (40), we apply the Lagrange duality method and derive  $f(\tau, \delta, \xi)$  by solving Problem (42). We note that with given dual variables  $\tau$ ,  $\delta$ , and  $\xi$ , Problem (42) can be decomposed into  $N_{as} + 1$  independent subproblems where one subproblem is for optimizing  $T$  and the other  $N_{as}$  subproblems are for

**Algorithm 3** Computing Capacity Allocation Algorithm for Solving Problem (40)

- 
- 1: Initialize  $\tau$ ,  $\delta$ , and  $\xi$ .
  - 2: **repeat**
  - 3:   Obtain the optimal  $T$  and  $f$  by solving (45) and (47), respectively;
  - 4:   Update the Lagrangian multipliers  $\tau$ ,  $\delta$ , and  $\xi$  by solving (48a), (48b) and (48c), respectively;
  - 5: **until** The objective function in (40a) converges.
- 

optimizing  $f_i, \forall i \in \mathcal{N}_{as}$ . The subproblem for optimizing  $T$  can be formulated as

$$\min_T T + \sum_{i=1}^{N_{as}} \frac{\tau_i I_i F_i}{T - \varpi_i} \quad (44a)$$

$$\text{s.t. } T \geq \varpi_i, \forall i \in \mathcal{N}_{as}. \quad (44b)$$

By setting the first-order derivative of (44a) with respect to  $T$  to zero, we observe that the optimal  $T$  should satisfy

$$T = \left\{ T \mid \sum_{i=1}^{N_{as}} \frac{\tau_i I_i F_i}{(T - \varpi_i)^2} = 1, T \geq \varpi_i \right\}, \quad (45)$$

which can be solved by applying the bisection search method.

Moreover, the subproblem for optimizing  $f_i, \forall i \in \mathcal{N}_{as}$  can be formulated as

$$\min_{f_i} \delta \kappa f_i^3 + \xi f_i - \tau_i f_i \quad (46a)$$

$$\text{s.t. } f_i \geq 0, \forall i \in \mathcal{N}_{as}. \quad (46b)$$

By setting the first-order derivative of (46a) with respect to  $f_i$  to zero, we obtain the closed-form expression of the optimal computing capacity allocation as

$$f_i = \left[ \sqrt{\frac{\tau_i - \xi}{3\delta\kappa}} \right]^+. \quad (47)$$

The value of dual variables  $\tau$ ,  $\delta$ , and  $\xi$  can be determined by the sub-gradient method. The updating procedure can be given by

$$\tau_i = \left[ \tau_i + \phi \left( \frac{I_i F_i}{T - \varpi_i} - f_i \right) \right]^+, \quad \forall i \in \mathcal{N}_{as} \quad (48a)$$

$$\delta = \left[ \delta + \phi \left( \sum_{i=1}^{N_{as}} \kappa f_i^3 - p_{max} + p_u \right) \right]^+ \quad (48b)$$

$$\xi = \left[ \xi + \phi \left( \sum_{i=1}^{N_{as}} f_i - f_{max} \right) \right]^+, \quad (48c)$$

where  $\phi > 0$  is the step-size in each iteration.

We summarize the procedures for obtaining the optimal solution to Problem (40) in Algorithm 3. Similar to the complexity analysis in Subsection III-B, we note that the total complexity of Algorithm 3 is  $\mathcal{O}(L_5 L_6 N_{as}^2 \log_2(1/\epsilon))$ , where  $L_5$  is the number of iterations for outer layer in Algorithm 3 and  $L_6$  is the number of iterations via the dual method of solving Problem (40).

### E. Caching Policy Subproblem

For any given  $\mathbf{y}, \theta, \eta, \mathbf{f}$ , and  $\mathbf{a}$ , the caching policy of Problem (15) can be optimized by solving the following problem

$$\min_{\mathbf{c}, T} T \quad (49a)$$

$$\text{s.t. } \varrho_i + \frac{(1 - c_i)I_i}{r_i^b} \leq T, \quad \forall i \in \mathcal{N} \quad (49b)$$

$$\sum_{i=1}^N c_i I_i \leq c_{max} \quad (49c)$$

$$c_i = \{0, 1\}, \quad \forall i \in \mathcal{N}, \quad (49d)$$

where  $\varrho_i = a_i \left( \frac{O_i}{r_i} + \frac{I_i F_i}{f_i} \right) + (1 - a_i) \left( \frac{I_i F_i}{f_i^{local}} + \frac{I_i}{r_i} \right)$ . Due to the linearity of the objective function and all constraints, we note that Problem (49) is a binary linear programming.

We first analyze the ideal scenario where the U-BS has a sufficiently large storage capability, i.e.,  $c_{max} \geq \sum_{i=1}^N I_i$ . In this case, since a lower maximum latency  $T$  might be achieved with a higher  $c_i$  according to (49b), we can easily obtain that the optimal solution for Problem (49) is  $c_i = 1, \forall i \in \mathcal{N}$ , i.e., the U-BS has cached the input data requested by all VR users and all requested input data can be directly obtained from the cache container of the U-BS without the backhaul transmissions, which significantly reduces the latency by eliminating the backhaul latency for all VR users.

For the general scenario where  $c_{max} \leq \sum_{i=1}^N I_i$ , due to the limited storage capability of the U-BS, only specific data which is requested by the VR users with high latency consumption will be pre-cached so that the maximum latency can be reduced via caching. Thus, to minimize the maximum latency  $T$ , we first sort the users based on the descending order in terms of  $\varrho_i + \frac{I_i}{r_i^b}$ . Next, we consider the input data required by the user with a higher  $\varrho_i + \frac{I_i}{r_i^b}$  will be cached at the U-BS with higher priority until the caching constraint cannot be satisfied. To derive the closed-form solution, we define a new indicator set  $\mathcal{S} \triangleq \{s_1, s_2, \dots, s_N\}$  which is sorted in a descending order in terms of  $\varrho_i + \frac{I_i}{r_i^b}$ , i.e.,  $s_1 = \arg \max_{\forall i \in \mathcal{N}} \varrho_i + \frac{I_i}{r_i^b}$  and  $s_N = \arg \min_{\forall i \in \mathcal{N}} \varrho_i + \frac{I_i}{r_i^b}$ . We further define the set  $\mathcal{S}_0 \triangleq \{s_1, s_2, \dots, s_{m-1}\}$ ,  $m = \min \{m : \sum_{i=1}^m I_{s_i} > c_{max}\}$ . By following [28], a closed-form expression for the optimal solution of (49) is given as

$$c_i = \begin{cases} 0, & \text{if } c_{max} \leq \sum_{i=1}^N I_i \text{ and } i \notin \mathcal{S}_0 \\ 1, & \text{otherwise.} \end{cases} \quad (50)$$

We note that the complexity for the caching policy subproblem is upper bounded by  $\mathcal{O}(N)$  and the actual complexity may be much smaller than this upper bound since the proposed approach may terminate when the caching storage constraint cannot be satisfied.

### F. Computing Policy Subproblem

For any given  $\mathbf{y}, \theta, \eta, \mathbf{f}$ , and  $\mathbf{c}$ , the computing policy Problem (15) can be optimized by solving the following problem

$$\min_{\mathbf{a}, T} T \quad (51a)$$

$$\text{s.t. } a_i v_i + o_i \leq T, \quad \forall i \in \mathcal{N} \quad (51b)$$

$$\sum_{i=1}^N a_i f_i \leq f_{max} \quad (51c)$$

$$p_u + \sum_{i=1}^N a_i \kappa f_i^3 \leq p_{max} \quad (51d)$$

$$a_i = \{0, 1\}, \quad \forall i \in \mathcal{N}, \quad (51e)$$

where  $v_i = \frac{O_i}{r_i} + \frac{I_i F_i}{f_i} - \frac{I_i F_i}{f_i^{local}} - \frac{I_i}{r_i} + \frac{I_i}{r_i^b}$  and  $o_i = \frac{I_i F_i}{f_i^{local}} + \frac{I_i}{r_i} + \frac{(1-c_i)I_i}{r_i^b}$ . We note that Problem (51) is a binary linear programming since the objective function and all constraints are linear.

To solve Problem (51), we first analyze the scenario when  $v_i \geq 0$ . In this case, we observe that a larger  $a_i$  might result in a higher  $T$  according to (51b). Thus, to minimize  $T$ , we can easily derive that  $a_i = 0$  when  $v_i \geq 0$ . This corresponds to the scenario that when the transmission and computing latencies at the  $i$ -th user of local computing is less than that of U-BS processing, the VR user chooses to self-execute its input data to reduce latency.

When  $v_i \leq 0$ , a lower maximum latency consumption might be achieved with a larger  $a_i$  according to (51b). However, due to the computing capacity and power constraints, only specific input data which is requested by the VR users with higher latency consumption will be processed at the U-BS so that the maximum latency can be minimized benefiting from the higher computing capacity at the U-BS. Thus, to minimize the maximum latency  $T$ , we first sort the users based on the descending order in terms of  $o_i$ . Next, we consider the input data of the user with a higher  $o_i$  will be processed at the U-BS with higher priority until the computing capacity constraint or the power constraint cannot be satisfied. To derive the closed-form solution, we define a new indicator set  $\mathcal{K} \triangleq \{k_1, k_2, \dots, k_N\}$  which is sorted in a descending order in terms of  $o_i$ , i.e.,  $k_1 = \arg \max_{\forall i \in \mathcal{N}} o_i$  and  $k_N = \arg \min_{\forall i \in \mathcal{N}} o_i$ . We further define the set  $\mathcal{K}_0 \triangleq \{k_1, k_2, \dots, k_{l-1}\}$ ,  $l = \min \{l_1, l_2\}$  where  $l_1 = \min \{l_1 : \sum_{i=1}^{l_1} f_{k_i} > f_{max}\}$  and  $l_2 = \min \{l_2 : p_u + \sum_{i=1}^{l_2} \kappa f_{k_i}^3 > p_{max}\}$ . Similar to the caching policy subproblem, a closed-form optimal solution of Problem (51) with complexity of  $\mathcal{O}(N)$  is given as

$$a_i = \begin{cases} 1, & \text{if } v_i \leq 0 \text{ and } i \in \mathcal{K}_0 \\ 0, & \text{otherwise.} \end{cases} \quad (52)$$

### G. Proposed Iterative Algorithm

The iterative procedure for solving Problem (15) is summarized in Algorithm 4, where the U-BS location, fronthaul and backhaul bandwidth allocation, computing capacity allocation, data caching policy and computing policy are successively optimized while keeping the other variables fixed until convergence, and the suboptimal solutions to Problem (15) can be obtained. In addition, the derived solution in each iteration will be applied as the input for the next iteration. We note that for U-BS location subproblem, since we only solve the approximated subproblem optimally, the convergence analysis for this subproblem should be studied.

Denote  $T_{UBS}^{ub}(\mathbf{y}[m], \theta[m], \eta[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m])$  as the objective values of (16). First, in step 3 of Algorithm 4, since the first-order Taylor expansions in (18) and (19) are tight

**Algorithm 4** Proposed Iterative Optimization for Problem (15)

- 1: Initialize  $m = 0$ ,  $\mathbf{y}[m]$ ,  $\boldsymbol{\theta}[m]$ ,  $\boldsymbol{\eta}[m]$ ,  $\mathbf{f}[m]$ ,  $\mathbf{c}[m]$ ,  $\mathbf{a}[m]$ .
- 2: **repeat**
- 3:   Given  $\{\boldsymbol{\theta}[m], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]\}$ , find the optimal U-BS location  $\mathbf{y}[m+1]$  by solving (20);
- 4:   Given  $\{\mathbf{y}[m+1], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]\}$ , find the optimal fronthaul bandwidth allocation  $\boldsymbol{\theta}[m+1]$  according to Algorithm 1;
- 5:   Given  $\{\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]\}$ , find the optimal backhaul bandwidth allocation  $\boldsymbol{\eta}[m+1]$  according to Algorithm 2;
- 6:   Given  $\{\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \boldsymbol{\eta}[m+1], \mathbf{c}[m], \mathbf{a}[m]\}$ , find the optimal computing capacity allocation  $\mathbf{f}[m+1]$  according to Algorithm 3;
- 7:   Given  $\{\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \boldsymbol{\eta}[m+1], \mathbf{f}[m+1], \mathbf{a}[m]\}$ , find the optimal caching policy  $\mathbf{c}[m+1]$  according to (50);
- 8:   Given  $\{\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \boldsymbol{\eta}[m+1], \mathbf{f}[m+1], \mathbf{c}[m+1]\}$ , find the optimal computing policy  $\mathbf{a}[m+1]$  according to (52);
- 9:   Update  $m = m + 1$ ;
- 10: **until** convergence.

bounds at given local point  $\mathbf{y}[m]$  for the original U-BS location subproblem (16), we have

$$\begin{aligned} & T(\mathbf{y}[m], \boldsymbol{\theta}[m], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]) \\ &= T_{UBS}^{ub}(\mathbf{y}[m], \boldsymbol{\theta}[m], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]). \end{aligned} \quad (53)$$

Notice that the U-BS location solution  $\mathbf{y}[m+1]$  for Problem (20) is optimal with other variables fixed, then it follows that

$$\begin{aligned} & T_{UBS}^{ub}(\mathbf{y}[m], \boldsymbol{\theta}[m], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]) \\ &\geq T_{UBS}^{ub}(\mathbf{y}[m+1], \boldsymbol{\theta}[m], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]) \\ &\geq T(\mathbf{y}[m+1], \boldsymbol{\theta}[m], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]), \end{aligned} \quad (54)$$

where the last inequality holds due to the fact that the objective value of Problem (20) is the upper bound of that of its original problem (16). Next, in step 4-8, since we solve the fronthaul and backhaul bandwidth allocation, computing capacity allocation, data caching policy and computing policy optimally, we have

$$\begin{aligned} & T(\mathbf{y}[m+1], \boldsymbol{\theta}[m], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]) \\ &\geq T(\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]) \\ &\geq T(\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \boldsymbol{\eta}[m+1], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]) \\ &\geq T(\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \boldsymbol{\eta}[m+1], \mathbf{f}[m+1], \mathbf{c}[m], \mathbf{a}[m]) \\ &\geq T(\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \boldsymbol{\eta}[m+1], \\ &\quad \mathbf{f}[m+1], \mathbf{c}[m+1], \mathbf{a}[m]) \\ &\geq T(\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \boldsymbol{\eta}[m+1], \\ &\quad \mathbf{f}[m+1], \mathbf{c}[m+1], \mathbf{a}[m+1]). \end{aligned} \quad (55)$$

According to (53)-(55), we can conclude that

$$\begin{aligned} & T(\mathbf{y}[m], \boldsymbol{\theta}[m], \boldsymbol{\eta}[m], \mathbf{f}[m], \mathbf{c}[m], \mathbf{a}[m]) \\ &\geq T(\mathbf{y}[m+1], \boldsymbol{\theta}[m+1], \boldsymbol{\eta}[m+1], \\ &\quad \mathbf{f}[m+1], \mathbf{c}[m+1], \mathbf{a}[m+1]), \end{aligned} \quad (56)$$

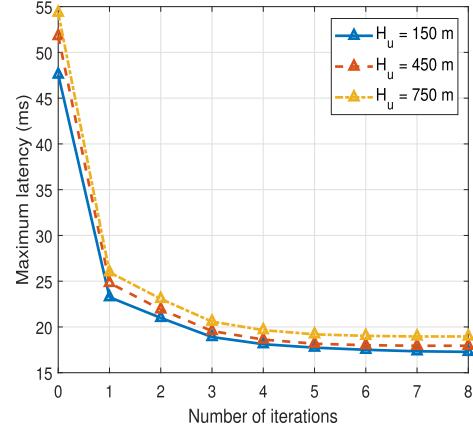


Fig. 2. Maximum latency versus number of iterations with different U-BS altitude  $H_u$ .

which shows that the algorithm yields a non-increasing sequence of the objective value. In addition, the objective value is lower bounded by zero. Hence, our proposed algorithm is guaranteed to converge. Although the obtained solution is generally suboptimal, we validate the effectiveness of our proposed Algorithm 4 in reducing the latency consumption via simulation results by comparing it with other benchmark strategies in Section IV. We note that the complexity of Algorithm 4 is the addition of the complexity in each step [20]. According to the aforementioned complexity analysis of each subproblem, we obtain that the overall complexity of Algorithm 4 is  $\mathcal{O}(L_1 L_2 N^2 \log_2(1/\epsilon) + L_3 L_4 N^2 \log_2(1/\epsilon) + L_5 L_6 N_{as}^2 \log_2(1/\epsilon) + 2N)$ , which shows that the complexity of Algorithm 4 is polynomial in the worst scenario.

#### IV. SIMULATION RESULTS

In this section, numerical results are presented to evaluate the performance of our proposed algorithm. We consider  $N = 6$  VR users that are randomly and uniformly distributed within a  $400\text{m} \times 400\text{m}$  square area. We set the altitude of U-BS as  $H_u = 150\text{ m}$ . The channel power gain is set as  $\beta_0 = 10^{-5}$ . We set the effective switched capacitance at the UAV as  $\kappa = 10^{-27}$  [21]. The noise spectral density is  $\sigma^2 = -169\text{ dBm/Hz}$ . The transmit powers at the U-BS and the cloud server are  $p_u = p_c = 0.5\text{ W}$ . We consider the input data size  $I_i$  follows a uniform distribution with  $I_i \sim U[10, 15]\text{ KB}$ , the ratio between  $O_i$  and  $I_i$  is set as  $\alpha = 2$ , and the required number of CPU cycles per bit is distributed as  $F_i \sim U[500, 800]\text{ cycles/bit}$ . The computing capacity of VR users follows a distribution of  $f_i^{local} \sim U[0.5, 1]\text{ GHz}$ . The maximum computing capacity, power budgets and caching storage of U-BS are set as  $f_{max} = 5\text{ GHz}$ ,  $P_{max} = 4\text{ W}$ , and  $c_{max} = 60\text{ KB}$ , respectively. The fronthaul and backhaul bandwidth are  $B = B_{back} = 1\text{ MHz}$ .

Fig. 2 shows the convergence behavior of Algorithm 4 with different U-BS altitude  $H_u$ . This figure shows that our proposed algorithm quickly converges within 8 iterations. Moreover, we observe that compared to its initial value, the maximum latency reduces by 63.7% from 47.6 ms to 17.3 ms when  $H_u = 150\text{ m}$ , which verifies the effectiveness of our proposed solution.

Fig. 3 shows the initial latency and optimized latency of each VR user where the initial latency is generated based on

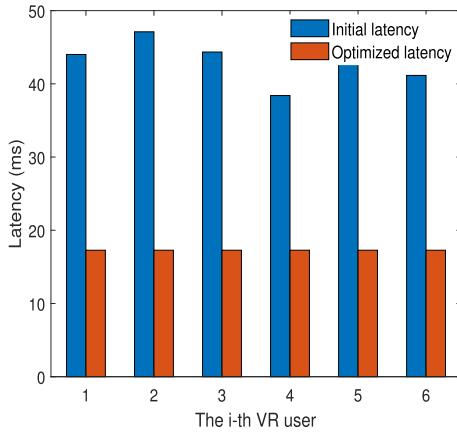
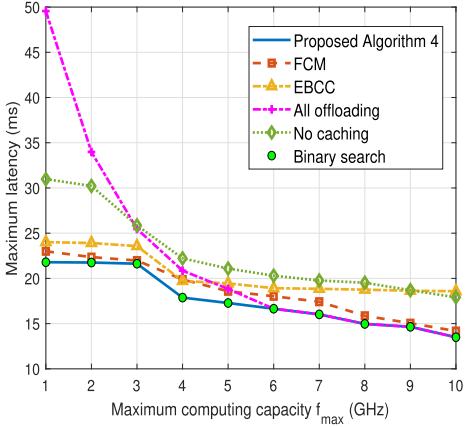
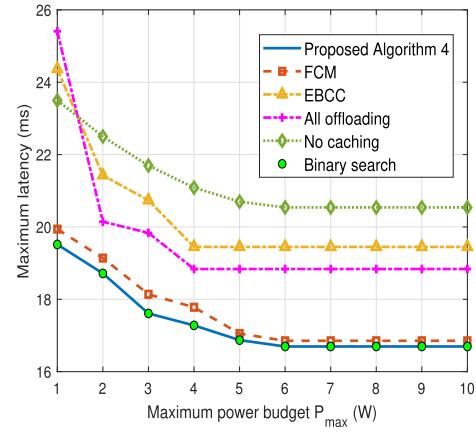


Fig. 3. The initial latency and optimized latency of each VR user.

Fig. 4. Maximum latency as a function of maximum computing capacity  $f_{max}$ .

one set of random realization of input data, computing capacity of each VR user and required number of CPU cycles per bit. It can be seen that our proposed joint optimization solution significantly reduces the maximum latency consumption among all VR users by comparing the initial latency and optimized latency. Moreover, we can see that the optimized latency of each VR user is almost equal, which shows that minimizing the maximum latency among all VR users is equivalent to guaranteeing the fairness among all VR users, so that the minimum quality-of-service can be improved.

In Fig. 4, we plot the maximum latency as a function of maximum computing capacity  $f_{max}$ . We compare our proposed Algorithm 4 with the following five benchmark schemes: 1) Fuzzy C-Means Clustering algorithm (FCM) [20]: the U-BS location is optimized based on FCM algorithm and all the other variables are optimized by using Algorithm 4; 2) Equal bandwidth and computing capacity (EBCC): We set  $\theta_i = 1/N$ ,  $f_i = a(i) * \min(f_{max}/N_{as}, ((P_{max} - p_u)/N_{as}/\kappa)^{1/3})$ ,  $\eta_i = (1 - c(i))/N_{uncached}$ ,  $\forall i \in \mathcal{N}$  and all the other variables are optimized by using Algorithm 4; 3) No caching: We set  $c_{max} = 0$  KB and all the other variables are optimized by using Algorithm 4; 4) All offloading: We set  $a_i = 1$ ,  $\forall i \in \mathcal{N}$  and all the other variables are optimized by using Algorithm 4; 5) Binary search: We solve the caching policy and computing policy subproblems by applying the binary search method and all

Fig. 5. Maximum latency as a function of maximum power budget  $P_{max}$ .

the other variables are optimized by using Algorithm 4. It can be seen that compared to the “Binary search” scheme which has an exponential complexity of  $\mathcal{O}(L_1 L_2 N^2 \log_2(1/\epsilon) + L_3 L_4 N^2 \log_2(1/\epsilon) + L_5 L_6 N_{as}^2 \log_2(1/\epsilon) + 2^{N+1})$ , our proposed Algorithm 4 with polynomial complexity achieves the same latency performance, which indicates that our proposed algorithm is stable and computationally efficient. Moreover, Fig. 4 shows that our proposed Algorithm 4 achieves a lower latency compared to other benchmark schemes except the “Binary search” scheme. Interestingly, we find that the performance gap between Algorithm 4 and “All offloading” scheme is significant when  $f_{max}$  is low, while it reduces to 0 when  $f_{max}$  is greater than 6 GHz. This is because when  $f_{max}$  is limited, e.g.,  $f_{max} = 1$  GHz, each VR user chooses to project the input data locally to reduce latency and the maximum latency is dominated by the computing latency which occupies 87%. While when  $f_{max} \geq 6$  GHz, the input data requested by all VR users will be processed at the U-BS, resulting in the same performance as “All offloading” scheme.

Fig. 5 shows the maximum latency as a function of maximum power budget  $P_{max}$ . It can be seen that our proposed Algorithm 4 achieves the same latency performance compared to the “Binary search” scheme and outperforms all the other baseline solutions. Moreover, we find that the maximum latency first decreases then keeps unchanged with an increasing  $P_{max}$ . This is because when  $P_{max}$  is limited, increasing  $P_{max}$  increases the computing resource allocated to process the input data requested by offloading VR users, resulting in a lower computing latency. However, when  $P_{max}$  is sufficient, the computing resource allocation is bounded by the maximum computing capacity, which makes the latency unchanged. In addition, we observe that caching is helpful to reduce the latency.

In Fig. 6, we plot the maximum latency as a function of fronthaul bandwidth  $B$ . We observe that our proposed Algorithm 4 achieves the same latency performance compared to the “Binary search” scheme and outperforms all the other baseline solutions. Moreover, we observe that when  $B$  is limited, e.g.,  $B = 0.5$  MHz, the maximum latency can be up to 26.5 ms and it mainly comes from the transmission latency, which occupies 57%. While when  $B = 1$  MHz, the maximum

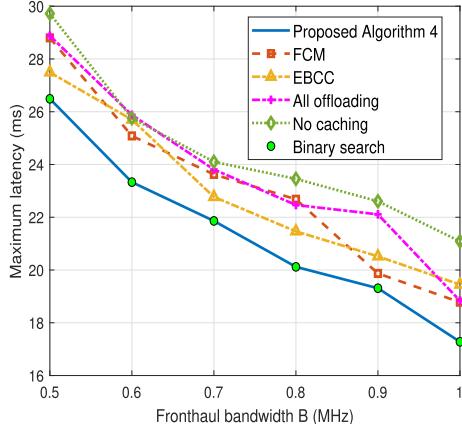


Fig. 6. Maximum latency as a function of fronthaul bandwidth  $B$ .

latency is 17.3 ms and the portion of transmission latency reduces to 27%.

## V. CONCLUSION

In this paper, we have presented the maximum latency minimization problem for a UAV-enabled communication, computing and caching VR delivery network. Specifically, we have jointly optimized the the U-BS location, fronthaul and backhaul bandwidth allocation, computing capacity allocation, caching and computing policies. To solve this nonconvex optimization problem, we have proposed an efficient iterative algorithm by applying the block coordinate descent method, the successive convex approximation technique and Lagrangian dual decomposition method. Simulation results demonstrated that our proposed algorithm significantly reduces the latency compared to benchmark schemes. Moreover, it can be seen that the maximum latency is mainly due to the transmission latency when the bandwidth is limited, whereas it is dominated by the computing latency when the computing resource is low. In addition, we showed that caching is helpful to reduce latency. We note that our work can be extended to consider that each VR user will execute multiple computation tasks to address the impact of queueing delay and consider the use of multiple UAVs to address limited battery capacity. Moreover, the extension to a more practical scenario that different VR users may require the same input data would be an interesting future research direction which results in a more-complex optimization problem with multiple possible computing and caching strategies.

## REFERENCES

- [1] M. Chen, W. Saad, C. Yin, and M. Debbah, "Data correlation-aware resource management in wireless virtual reality (VR): An echo state transfer learning approach," *IEEE Trans. Commun.*, vol. 67, no. 6, pp. 4267–4280, Jun. 2019.
- [2] M. Chen, W. Saad, and C. Yin, "Virtual reality over wireless networks: Quality-of-service model and learning-based resource management," *IEEE Trans. Commun.*, vol. 66, no. 11, pp. 5621–5635, Nov. 2018.
- [3] M. S. Elbamby, C. Perfecto, M. Bennis, and K. Doppler, "Toward low-latency and ultra-reliable virtual reality," *IEEE Netw.*, vol. 32, no. 2, pp. 78–84, Mar./Apr. 2018.
- [4] E. Bastug, M. Bennis, M. Médard, and M. Debbah, "Toward interconnected virtual reality: Opportunities, challenges, and enablers," *IEEE Commun. Mag.*, vol. 55, no. 6, pp. 110–117, Jun. 2017.
- [5] X. Yang *et al.*, "Communication-constrained mobile edge computing systems for wireless virtual reality: Scheduling and tradeoff," *IEEE Access*, vol. 6, pp. 16665–16677, 2018.
- [6] J. Park, P. Popovski, and O. Simeone, "Minimizing latency to support VR social interactions over wireless cellular systems via bandwidth allocation," *IEEE Wireless Commun. Lett.*, vol. 7, no. 5, pp. 776–779, Oct. 2018.
- [7] A. Al-Shuwaili and O. Simeone, "Energy-efficient resource allocation for mobile edge computing-based augmented reality applications," *IEEE Wireless Commun. Lett.*, vol. 6, no. 3, pp. 398–401, Jun. 2017.
- [8] Y. Dai, D. Xu, S. Maharjan, and Y. Zhang, "Joint computing offloading and user association in multi-task mobile edge computing," in *IEEE Trans. Veh. Technol.*, vol. 67, no. 12, pp. 12313–12325, Dec. 2018.
- [9] Z. Xiang, F. Gabriel, E. Urbano, G. T. Nguyen, M. Reisslein, and F. H. P. Fitzek, "Reducing latency in virtual machines: Enabling tactile Internet for human-machine co-working," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 5, pp. 1098–1116, May 2019.
- [10] P. Wang, Z. Zheng, B. Di, and L. Song, "HetMEC: Latency-optimal task assignment and resource allocation for heterogeneous multi-layer mobile edge computing," *IEEE Trans. Wireless Commun.*, vol. 18, no. 10, pp. 4942–4956, Oct. 2019.
- [11] Y. Zhou *et al.*, "Offloading optimization for low-latency secure mobile edge computing systems," *IEEE Wireless Commun. Lett.*, vol. 9, no. 4, pp. 480–484, Apr. 2020.
- [12] T. Dang and M. Peng, "Joint radio communication, caching, and computing design for mobile virtual reality delivery in fog radio access networks," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 7, pp. 1594–1607, Jul. 2019.
- [13] Y. Sun, Z. Chen, M. Tao, and H. Liu, "Communications, caching, and computing for mobile virtual reality: Modeling and tradeoff," *IEEE Trans. Commun.*, vol. 67, no. 11, pp. 7573–7586, Nov. 2019.
- [14] J. Zhang, X. Hu, Z. Ning, E. C.-H. Ngai, L. Zhou, J. Wei, J. Cheng, B. Hu, and V. C. M. Leung, "Joint resource allocation for latency-sensitive services over mobile edge computing networks with caching," *IEEE Internet Things J.*, vol. 6, no. 3, pp. 4283–4294, Jun. 2019.
- [15] Q. Wu, Y. Zeng, and R. Zhang, "Joint trajectory and communication design for multi-UAV enabled wireless networks," *IEEE Trans. Wireless Commun.*, vol. 17, no. 3, pp. 2109–2121, Mar. 2018.
- [16] Y. Zhou *et al.*, "Improving physical layer security via a UAV friendly jammer for unknown eavesdropper location," *IEEE Trans. Veh. Technol.*, vol. 67, no. 11, pp. 11280–11284, Nov. 2018.
- [17] C. Pan, H. Ren, Y. Deng, M. Elkashlan, and A. Nallanathan, "Joint blocklength and location optimization for URLLC-enabled UAV relay systems," *IEEE Commun. Lett.*, vol. 23, no. 3, pp. 498–501, Mar. 2019.
- [18] Z. Yang *et al.*, "Joint altitude, beamwidth, location, and bandwidth optimization for UAV-enabled communications," *IEEE Commun. Lett.*, vol. 22, no. 8, pp. 1716–1719, Aug. 2018.
- [19] Y. Zhou *et al.*, "Secure Communications for UAV-enabled mobile edge computing systems," *IEEE Trans. Commun.*, vol. 68, no. 1, pp. 376–388, Jan. 2020.
- [20] Z. Yang, C. Pan, K. Wang, and M. Shikh-Bahaei, "Energy efficient resource allocation in UAV-enabled mobile edge computing networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 9, pp. 4576–4589, Sep. 2019.
- [21] Q. Hu, Y. Cai, G. Yu, Z. Qin, M. Zhao, and G. Y. Li, "Joint offloading and trajectory design for UAV-enabled mobile edge computing systems," *IEEE Internet Things J.*, vol. 6, no. 2, pp. 1879–1892, Apr. 2019.
- [22] N. Zhao *et al.*, "Caching UAV assisted secure transmission in hyper-dense networks based on interference alignment," *IEEE Trans. Commun.*, vol. 66, no. 5, pp. 2281–2294, May 2018.
- [23] M. Chen, M. Mozaffari, W. Saad, C. Yin, M. Debbah, and C. S. Hong, "Caching in the sky: Proactive deployment of cache-enabled unmanned aerial vehicles for optimized quality-of-experience," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 5, pp. 1046–1061, May 2017.
- [24] X. Xu, Y. Zeng, Y. L. Guan, and R. Zhang, "Overcoming endurance issue: UAV-enabled communications with proactive caching," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 6, pp. 1231–1244, Jun. 2018.
- [25] Q. Wu and R. Zhang, "Common throughput maximization in UAV-enabled OFDMA systems with delay consideration," *IEEE Trans. Wireless Commun.*, vol. 66, no. 12, pp. 6614–6627, Dec. 2018.
- [26] Y. Zeng, X. Xu, and R. Zhang, "Trajectory design for completion time minimization in UAV-enabled multicasting," *IEEE Trans. Wireless Commun.*, vol. 17, no. 4, pp. 2233–2246, Apr. 2018.
- [27] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [28] J. Dai, Z. Hu, B. Li, J. Liu, and B. Li, "Collaborative hierarchical caching with dynamic request routing for massive content distribution," in *Proc. IEEE INFOCOM*, Dec. 2012, pp. 2444–2452.



**Yi Zhou** received the Ph.D. degree from the School of Engineering and Information Technologies, The University of Sydney, Australia, in 2020. Her research interests include physical layer security, UAV communications, and 5G related communications. She was a recipient of the Post-Graduate Scholarship and the Norman I. Price Scholarship from the Center of Excellence in Telecommunications, School of Electrical and Information Engineering, The University of Sydney.



**Cunhua Pan** (Member, IEEE) received the B.S. and Ph.D. degrees from the School of Information Science and Engineering, Southeast University, Nanjing, China, in 2010 and 2015, respectively.

From 2015 to 2016, he was a Research Associate with the University of Kent, U.K. He held a post-doctoral position with the Queen Mary University of London, U.K., from 2016 and 2019, where he is currently a Lecturer. His research interests mainly include intelligent reflection surface (IRS), ultra-reliable low latency communication (URLLC), machine learning, UAV, the Internet of Things, and mobile edge computing. He serves as a TPC member for numerous conferences, such as the ICC and GLOBECOM, and the Student Travel Grant Chair for ICC 2019. He also serves as an Editor of IEEE WIRELESS COMMUNICATION LETTERS, IEEE COMMUNICATIONS LETTERS, and IEEE ACCESS.



**Phee Lep Yeo** (Member, IEEE) received the B.E. degree (with University Medal) and the Ph.D. degree from The University of Sydney (USYD), Australia, in 2004 and 2012, respectively.

From 2008 to 2012, he was with the Telecommunications Laboratory, The University of Sydney, and the Wireless and Networking Technologies Laboratory at the Commonwealth Scientific and Industrial Research Organization (CSIRO), Australia. From 2012 to 2016, he was with the Department of Electrical and Electronic Engineering, The University of Melbourne, Australia. Since 2016, he has been a Senior Lecturer with the School of Electrical and Information Engineering, USYD. His current research interests include secure wireless communications, ultra-reliable and low-latency communications (URLLC), ultra-dense networks, and multiscale molecular communications. He is a recipient of the 2020 USYD Robinson Fellowship, the 2017 Alexander von Humboldt Research Fellowship for Experienced Researchers, and the 2014 Australian Research Council (ARC) Discovery Early Career Researcher Award (DECRA). He received Best Paper Awards at IEEE ICC 2014 and IEEE VTC-Spring 2013, and the Best Student Paper awards at AusCTW 2013 and 2019. He has served as the TPC Chair for the 2016 Australian Communications Theory Workshop (AusCTW) and a TPC member for IEEE GLOBECOM, ICC, and VTC conferences.



**Kezhi Wang** (Senior Member, IEEE) received the B.E. and M.E. degrees from the School of Automation, Chongqing University, China, in 2008 and 2011, respectively, and the Ph.D. degree in engineering from the University of Warwick, U.K., in 2015. He was a Senior Research Officer with the University of Essex, U.K. He is currently a Senior Lecturer with the Department of Computer and Information Sciences, Northumbria University, U.K. His research interests include wireless communication, mobile edge computing, UAV communication, and machine learning.



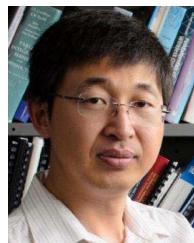
**Maged Elkashlan** (Senior Member, IEEE) received the Ph.D. degree in electrical engineering from The University of British Columbia, Canada, 2006.

From 2007 to 2011, he was with the Commonwealth Scientific and Industrial Research Organization (CSIRO), Australia. During this time, he held visiting appointments at the University of New South Wales and University of Technology Sydney. In 2011, he joined the School of Electronic Engineering and Computer Science, Queen Mary University of London, U.K. His research interests

fall into the broad areas of communication theory and statistical signal processing. He received the Best Paper Awards at IEEE International Conference on Communications (ICC) in 2016 and 2014, the International Conference on Communications and Networking in China (CHINACOM) in 2014, and IEEE Vehicular Technology Conference (VTC-Spring) in 2013. He is an Editor of IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY and IEEE TRANSACTIONS ON MOLECULAR, BIOLOGICAL AND MULTI-SCALE COMMUNICATIONS.



**Branka Vucetic** (Life Fellow, IEEE) is currently an ARC Laureate Fellow and the Director of the Centre of Excellence for IoT and Telecommunications, The University of Sydney. Her current work is in the areas of wireless networks and the Internet of Things. In the area of wireless networks, she works on communication system design for millimeter-wave (mmWave) frequency bands. In the area of the Internet of things, she works on providing wireless connectivity for mission-critical applications. She is a fellow of the Australian Academy of Science, the Australian Academy of Technological Sciences and Engineering, and the Engineers Australia.



**Yonghui Li** (Fellow, IEEE) received the Ph.D. degree from the Beijing University of Aeronautics and Astronautics in November 2002.

From 1999 to 2003, he was affiliated with Linkair Communication Inc., where he held the position of Project Manager with responsibility for the design of physical layer solutions for the LAS-CDMA system. Since 2003, he has been with the Centre of Excellence in Telecommunications, The University of Sydney, Australia. He is currently a Professor with the School of Electrical and Information Engineering, University of Sydney. He is a recipient of the Australian Queen Elizabeth II Fellowship in 2008 and the Australian Future Fellowship in 2012. His current research interests are in the area of wireless communications, with a particular focus on MIMO, millimeter-wave communications, machine to machine communications, coding techniques, and cooperative communications. He holds a number of patents granted and pending in these fields. He is currently an Editor for IEEE TRANSACTIONS ON COMMUNICATIONS and IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY. He was also the Guest Editor for IEEE JSAC Special Issue on Millimeter Wave Communications for Future Mobile Networks. He received the Best Paper Awards from IEEE International Conference on Communications (ICC) 2014, IEEE PIMRC 2017, and IEEE Wireless Days Conferences (WD) 2014.