# Reconfigurable Intelligent Surface-assisted Edge Computing to Minimize Delay in Task Offloading

Mithun Mukherjee*, Vikas Kumar†, Suman Kumar‡, Jaime Lloret§, Qi Zhang¶, and Mian Guo‖

*School of Artificial Intelligence, Nanjing University of Information Science and Technology, China, m.mukherjee@ieee.org
†Bharat Sanchar Nigam Limited, India, vikas.kr@bsnl.co.in
‡Department of Mathematics, IGNTU Amarkantak, MP, India, suman@igntu.ac.in
§Universitat Politecnica de Valencia, Spain, jlloret@dcom.upv.es
¶DIGIT, Department of Electrical and Computer Engineering, Aarhus University, Aarhus, Denmark, qz@ece.au.dk
‖School of Electronics and Information, Guangdong Polytechnic Normal University, P.R. China, mian.guo@ieee.org

*Abstract*—The advantage of computational resources in edge computing near the data source has kindled growing interest in delay-sensitive Internet of Things (IoT) applications. However, the benefit of the edge server is limited by the uploading and downloading links between end-users and edge servers when these end-users seek computational resources from edge servers. The scenario becomes more severe when the user-end's devices are in the shaded region resulting in low uplink/downlink quality. In this paper, we consider a reconfigurable intelligent surface (RIS)-assisted edge computing system, where the benefits of RIS are exploited to improve the uploading transmission rate. We further aim to minimize the delay of worst-case in the network when the end-users either compute task data in their local CPU or offload task data to the edge server. Next, we optimize the uploading bandwidth allocation for every end-user's task data to minimize the maximum delay in the network. The above optimization problem is formulated as quadratically constrained quadratic programming. Afterward, we solve this problem by semidefinite relaxation. Finally, the simulation results demonstrate that the proposed strategy is scalable under various network settings.

## I. INTRODUCTION

In recent years, several industries have been focusing their technological advancement towards high performance computing in cloud data centers. For example, in 2020, NVDIA announced the potential use of DGX A100 NVIDIA's third-generation Artificial Intelligence (AI) system box [1] that is aimed at the massive gain in performance for AI-related and cutting-edge applications with less power consumption. At the same time, we are witnessing the paradigm changing from constituting a well-run centralized data center infrastructure to the network edge [2]–[5], particularly when there is a need to deliver proximity, low-latency, and reliable services for the mission-critical applications, such as remote-surgery, industrial automation and driverless cars. The leading industries with their cloud service providers (e.g., EGX Edge AI platform, NVIDIA RTX graphics with CloudXR, GPU virtualization, and Qualcomm Technologies' Boundless XR client optimizations [6] and EdgeConneX [7]) are making their way for the deployment of edge-assisted service provisioning.

### A. Motivation

Although MEC brings computational, caching, and storage resources towards the network edge, the connectivity and coverage of the access points and base station play critical roles. To say, when end-users aim to avail the computing, caching, and storage resources of edge server, they need to rely on the wireless channels. Basically, irrespective of application type and service, the uploading/downloading time is an important factor in delay-sensitive service provisioning. This becomes even worse when the network coverage is poor near the cell edge or blocked by obstacles. The uploading and downloading rate and the resulting latency are significantly affected by communication resource allocation. This, in turn, affects the computation delay of end-user's task. To address the above shortcoming that arises due to the connectivity and coverage of the network services, reconfigurable intelligent surface (RIS) [8] can assist MEC.

### B. Related Work

The role of RIS has been studied in the MEC system, where the end-users aim to offload their computation-intensive tasks to the edge server that resides at the access point [9], [10]. They formulated a latency minimization problem by optimizing the task offloading data size, edge computing resource allocation and RIS phase shift coefficients. To maximize the total amount of data (in terms of bits) processed by end-users and edge server, Chu et al. [11] suggested how to adjust the phase shift of the RIS in addition to the transmit power and time allocation for the end-users, and edge server's computing resource allocation for the end-users. Another study in [12] shows how the edge server adjusts the RIS controller to maximize its revenue while guaranteeing the customized information rate for each end-user. Later, another parallel work studied the RIS-enabled MEC system in [13]. Again, this was to minimize the latency, which is basically calculated as the sum of two end-user's computation offloading time. Moreover, Cao et al. [14] have shown how RIS can resolve the link blockage problem in the mm-Wave MEC system to guarantee real-time offloading from the end-users. This is an interesting and detailed study on how RIS can directly affect the task offloading chances for the end-users that suffer from mm-Wave link blockage.

Recently, over-the-air computation (AirComp) [15] that integrates communication and computation has attracted
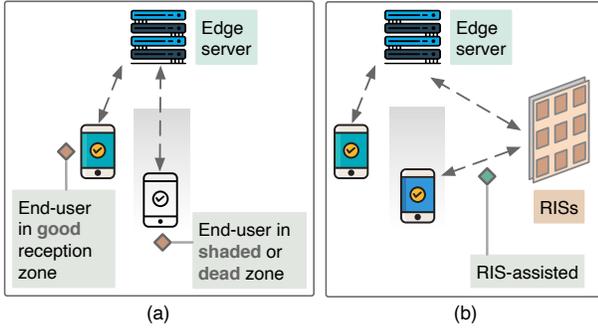
Fig. 1. (a) An illustration of a RIS-assisted edge computing system with end-users under *good* and *shadowed* region. (b) The RIS can assist to improve the uplink and downlink quality for the end-users under shadowed region.

academia and industries' attention due to its fast data aggregation from IoT devices. However, due to the unreliable channel conditions, the performance of AirComp is severely limited. To address this, RIS [16]–[18] has been found a suitable candidate to assist the uplink and downlink transmission.

### C. Our Contributions and Organization

We summarize our main contributions as follows: We consider an RIS-assisted edge computing system, where end-users offload their task data to the edge server to minimize the overall delay. We aim to leverage the benefits of RIS for the uplink transmission rate in data offloading to the edge server. With the assistance of RIS, a delay-minimization problem is formulated by optimizing the offloading decision variables and bandwidth allocation for the offloaded task data. We formulate the above optimization problem as Quadratically Constrained Quadratic Programming (QCQP) problem. Afterward, we apply semi-definite relaxation (SDR) to solve the problem. Finally, we show that the proposed offloading strategy with RIS can achieve better performance than without RIS assistance and *local CPU only* approaches.

The rest of the paper is organized as follows. In Section II, we discuss the RIS-assisted MEC system model. We formulate the optimization problem and apply SDR in Section III. The simulation results are presented in Section IV. Finally, we conclude our work in Section V.

## II. SYSTEM MODEL

We consider an RIS-assisted edge computing system, as shown in Fig. 1. The set of the end-users is denoted as $\mathcal{M} = \{1, 2, \ldots, M\}$, where $M$ is the total number of end-users in the network. Due to the limited computational resources (in terms of CPU speed) in local CPU, these end-users often offload task data to edge server when the tasks demand fast processing. Note that among these end-users, we assume that $K$ end-users have good reception quality and the remaining $(M - K)$ end-users are in poor signal reception area. We denote the end-users in *good* and *poor* signal reception areas as the $i$th and the $j$th end-user, respectively, where $i = 1, \ldots, K$ and $j = (K + 1), \ldots, M$. We further denote the offloading decision variable for the $m$th end-user at local device of each

end-user as $x_m$ and the task processing decision variable at edge server for the end-user as $y_m$, where $m = 1, \ldots, M$ and

$$
x_m = \begin{cases} 1 & \text{when the } m\text{th end-user's task data is} \\ & \text{locally processed,} \\ 0 & \text{otherwise,} \end{cases}
$$

$$
y_m = \begin{cases} 1 & \text{when the } m\text{th end-user's task data is} \\ & \text{offloaded and processed at the edge server,} \\ 0 & \text{otherwise.} \end{cases}
$$

Note that $x_m + y_m = 1$. Moreover, these binary decision variables satisfy $x_m(1 - x_m) = 0$ and $y_m(1 - y_m) = 0$.

### A. Local Computing Delay

When a task is locally processed by the end-user, the computation delay becomes

$$
T_i^L = \frac{x_i D_i L}{f_i^l} \text{ [s]} \quad i \in \{1, 2, \ldots, K\}, \tag{1a}
$$

$$
T_j^L = \frac{x_j D_j L}{f_j^l} \text{ [s]} \quad j \in \{(K+1), \ldots, M\}, \tag{1b}
$$

where $D_i$ and $D_j$ is the input data size [bits] of the $i$th and $j$th end-user, respectively, $L$ is the processing density [CPU cycles/bit] for a task, and $f_i$ and $f_j$ denote the CPU clock speed [CPU cycles/s] of the $i$th and $j$th end-user, respectively. We assume equal task processing density for every end-users.

### B. Offloading Delay for End-users without RIS Assistance

Basically, this is the case when the end-users are in *good* signal reception area. When the $i$th end-user task data is offloaded and processed at the edge server, the offloading delay becomes

$$
T_i^{\mathsf{E}} = y_i \left( \underbrace{\frac{D_i}{\eta_i \beta_i C}}_{\text{uploading}} + \underbrace{\frac{D_i L}{f_i^e}}_{\text{computation}} \right) \text{ [s]}, \tag{2}
$$

where $\eta_i$ is the spectral efficiency of uplink transmission between the $i$th end-user and edge server, $C$ is the total uplink bandwidth, $\beta_i$ is the fraction of total uplink bandwidth allocated to the $i$th end-user, and $f_i^e$ is the CPU rate allocated by the edge server to process the $i$th end-user's offloaded task.

### C. Offloading Delay for End-users with RIS Assistance

When the $j$th end-user with poor wireless connection offloads its task data to the edge server with RIS assistance, the offloading delay can be written as

$$
T_j^{\mathsf{E}} = y_j \left( \underbrace{\frac{D_j}{\eta_j \beta_j C}}_{\text{uploading}} + \underbrace{\frac{D_j L}{f_j^e}}_{\text{computation}} \right) \text{ [s]}, \tag{3}
$$

where $\eta_j$ is the spectral efficiency of RIS-assisted uplink transmission between the $j$th end-user and the edge server, $\beta_j$ is the fractional value of total uplink bandwidth allocated to the $j$th end-user and $f_j^e$ is the CPU rate allocated by edge server to process the $j$th end-user's offloaded task.

## III. PROBLEM FORMULATION

We write the delay of the worst case in the network as $\max\{(T_i^{\mathsf{L}} + T_i^{\mathsf{E}}), (T_j^{\mathsf{L}} + T_j^{\mathsf{E}})\}\forall i \in \mathcal{N}_1, j \in \mathcal{N}_2$, where $\mathcal{N}_1 = \{1\ldots K\}$, $\mathcal{N}_2 = \{K+1, \ldots M\}$. We aim to minimize the maximum delay by jointly optimizing the task offloading decision vector $\boldsymbol{\xi} = [x_m, y_m]^{\mathsf{T}}$ and the bandwidth allocation vector $\mathbf{r} = [\beta_m]^{\mathsf{T}}$, where

$$T_i^{\mathsf{L}} + T_i^{\mathsf{E}} = \frac{D_i\,L\,x_i}{f_i^l} + \frac{D_i\,y_i}{\eta_i\,\beta_i\,C} + \frac{D_i\,L\,y_i}{f_i^e}\ [\text{s}]\,, \quad (4a)$$

$$T_j^{\mathsf{L}} + T_j^{\mathsf{E}} = \frac{D_j L x_j}{f_j^l} + \frac{D_j\,y_j}{\eta_j\,\beta_j\,C} + \frac{D_j L y_j}{f_j^e}\ [\text{s}]\,. \quad (4b)$$

We define the above optimization problem as

$$\min_{\boldsymbol{\xi}, \mathbf{r}} \quad \max\{(T_i^{\mathsf{L}} + T_i^{\mathsf{E}}), (T_j^{\mathsf{L}} + T_j^{\mathsf{E}})\} \ \forall\ i \in \mathcal{N}_1, j \in \mathcal{N}_2$$
$$(5a)$$

$$\text{s.t.}\quad x_m(1 - x_m) = 0, \quad (5b)$$
$$y_m(1 - y_m) = 0, \quad (5c)$$
$$x_m + y_m = 1, \quad (5d)$$
$$\sum_{i=1}^{K} \beta_i + \sum_{j=K+1}^{M} \beta_j \leq 1, \quad (5e)$$

where the constraint (5e) corresponds to the total uplink bandwidth $C$. Now, we take an auxiliary variable $t$ as

$$\max_{i \in \mathcal{N}_1, j \in \mathcal{N}_2} \left\{(T_i^{\mathsf{L}} + T_i^{\mathsf{E}}), (T_j^{\mathsf{L}} + T_j^{\mathsf{E}})\right\} = t\,, \quad (6)$$

then, from (4) and (6), we write

$$\frac{D_i\,L\,x_i\,\beta_i}{f_i^l} + \frac{D_i\,y_i}{\eta_i\,C} + \frac{D_i\,L\,y_i\,\beta_i}{f_i^e} - \beta_i\,t \leq 0, \quad (7a)$$

$$\frac{D_j L x_j \beta_j}{f_j^l} + \frac{D_j y_j}{\eta_j\,C} + \frac{D_j L y_j \beta_j}{f_j^e} - \beta_j\,t \leq 0. \quad (7b)$$

Accordingly, the optimization problem becomes

$$\min_{\boldsymbol{\xi}, \mathbf{r}} \quad t \quad (8a)$$
$$\text{s.t.}\quad x_m(1 - x_m) = 0, \quad (8b)$$
$$y_m(1 - y_m) = 0, \quad (8c)$$
$$x_m + y_m = 1, \quad (8d)$$
$$\sum_{i=1}^{K} \beta_i + \sum_{j=K+1}^{M} \beta_j \leq 1, \quad (8e)$$
$$(7a) \text{ and } (7b)\,. \quad (8f)$$

### A. Vector-matrix Formation

Now, we denote $\mathbf{w} = [x_1, x_2, \ldots, x_K, x_{K+1}, \ldots, x_M, y_1, y_2, \ldots, y_K, y_{K+1}, \ldots, y_M, \beta_1, \beta_2, \ldots, \beta_K, \beta_{K+1}, \ldots, \beta_M, t]^{\mathsf{T}}$ and define the unit vector as $\mathbf{e}_q = [\mathbf{0}_{1\times(q-1)}, 1,$

$\mathbf{0}_{1\times(3M+1-q)}]^{\mathsf{T}}$. Then, the matrix form of problem (8) can be expressed as

$$\min_{\mathbf{w}} \quad \mathbf{e}_{(3M+1)}^{\mathsf{T}}\mathbf{w} \quad (9a)$$
$$\text{s.t.}\quad \mathbf{w}^{\mathsf{T}}\mathbf{A}_{x,m}\mathbf{w} - \mathbf{e}_m^{\mathsf{T}}\mathbf{w} = 0, \quad (9b)$$
$$\mathbf{w}^{\mathsf{T}}\mathbf{A}_{y,m}\mathbf{w} - \mathbf{e}_{m+M}^{\mathsf{T}}\mathbf{w} = 0, \quad (9c)$$
$$\mathbf{e}_m^{\mathsf{T}}\mathbf{w} + \mathbf{e}_{m+M}^{\mathsf{T}}\mathbf{w} = 1, \quad (9d)$$
$$\sum_{i=1}^{K} \mathbf{e}_{i+2M}^{\mathsf{T}}\mathbf{w} + \sum_{j=K+1}^{M} \mathbf{e}_{j+2M}^{\mathsf{T}}\mathbf{w} \leq 1, \quad (9e)$$
$$\mathbf{w}^{\mathsf{T}}\mathbf{A}_{\beta x,i}\mathbf{w} + \mathbf{w}^{\mathsf{T}}\mathbf{A}_{\beta y,i}\mathbf{w} + \mathbf{b}_{cy,i}^{\mathsf{T}}\mathbf{w}$$
$$+ \mathbf{w}^{\mathsf{T}}\mathbf{A}_{\beta t,i}\mathbf{w} \leq 0\,, \quad (9f)$$
$$\mathbf{w}^{\mathsf{T}}\mathbf{A}_{\beta x,j}\mathbf{w} + \mathbf{w}^{\mathsf{T}}\mathbf{A}_{\beta y,j}\mathbf{w} + \mathbf{b}_{cy,j}^{\mathsf{T}}\mathbf{w}$$
$$+ \mathbf{w}^{\mathsf{T}}\mathbf{A}_{\beta t,j}\mathbf{w} \leq 0\,, \quad (9g)$$

where

$$\mathbf{A}_{x,m} = \begin{bmatrix} \mathbf{0}_{(m-1)\times(3M+1)} \\ \hline \mathbf{e}_m^{\mathsf{T}} \\ \hline \mathbf{0}_{(3M+1-m)\times(3M+1)} \end{bmatrix},$$

$$\mathbf{A}_{y,m} = \begin{bmatrix} \mathbf{0}_{(M-1+m)\times(3M+1)} \\ \hline \mathbf{e}_{(m+M)}^{\mathsf{T}} \\ \hline \mathbf{0}_{(2M+1-m)\times(3M+1)} \end{bmatrix},$$

$$\mathbf{b}_{cy,i} = k_i^c\mathbf{e}_{M+i}, \ k_i^l = \frac{D_i\,L}{f_i^l}, \ k_i^e = \frac{D_i\,L}{f_i^e}, \ k_i^c = \frac{D_i}{\eta_i C},$$

$$\mathbf{A}_{\beta x,i} = \frac{k_i^l}{2}\begin{bmatrix} \mathbf{0}_{(i-1)\times(3M+1)} \\ \hline \mathbf{e}_{i+2M}^{\mathsf{T}} \\ \hline \mathbf{0}_{(2M-1)\times(3M+1)} \\ \hline \mathbf{e}_i^{\mathsf{T}} \\ \hline \mathbf{0}_{(M+1-i)\times(3M+1)} \end{bmatrix},$$

$$\mathbf{A}_{\beta y,i} = \frac{k_i^e}{2}\begin{bmatrix} \mathbf{0}_{(M-1+i)\times(3M+1)} \\ \hline \mathbf{e}_{i+2M}^{\mathsf{T}} \\ \hline \mathbf{0}_{(M-1)\times(3M+1)} \\ \hline \mathbf{e}_{i+M}^{\mathsf{T}} \\ \hline \mathbf{0}_{(M+1-i)\times(3M+1)} \end{bmatrix},$$

$$\mathbf{b}_{cy,j} = k_j^c\mathbf{e}_{M+j}, \ k_j^l = \frac{D_j\,L}{f_j^l}, \ k_j^e = \frac{D_j\,L}{f_j^e}, \ k_j^c = \frac{D_j}{\eta_j C},$$

$$\mathbf{A}_{\beta t,i} = -\frac{1}{2}\begin{bmatrix} \mathbf{0}_{(2M-1+i)\times(3M+1)} \\ \hline \mathbf{e}_{3M+1}^{\mathsf{T}} \\ \hline \mathbf{0}_{(M-i)\times(3M+1)} \\ \hline \mathbf{e}_{i+2M}^{\mathsf{T}} \end{bmatrix},$$

$$\mathbf{A}_{\beta x,j} = \frac{k_j^l}{2}\begin{bmatrix} \mathbf{0}_{(j-1)\times(3M+1)} \\ \hline \mathbf{e}_{j+2M}^{\mathsf{T}} \\ \hline \mathbf{0}_{(2M-1)\times(3M+1)} \\ \hline \mathbf{e}_j^{\mathsf{T}} \\ \hline \mathbf{0}_{(M+1-j)\times(3M+1)} \end{bmatrix},$$

$$\mathbf{A}_{\beta y,j} = \frac{k_j^e}{2}\begin{bmatrix} \mathbf{0}_{(M-1+j)\times(3M+1)} \\ \hline \mathbf{e}_{j+2M}^{\mathsf{T}} \\ \hline \mathbf{0}_{(M-1)\times(3M+1)} \\ \hline \mathbf{e}_{j+M}^{\mathsf{T}} \\ \hline \mathbf{0}_{(M+1-j)\times(3M+1)} \end{bmatrix},$$

$$\mathbf{A}_{\beta t,j} = -\frac{1}{2} \begin{bmatrix} \mathbf{0}_{(2M-1+j)\times(3M+1)} \\ \hline \mathbf{e}_{3M+1}^{\mathsf{T}} \\ \hline \mathbf{0}_{(M-j)\times(3M+1)} \\ \hline \mathbf{e}_{j+2M}^{\mathsf{T}} \end{bmatrix}.$$

### B. QCQP Formulation

Defining $\mathbf{z} = [\mathbf{w}^{\mathsf{T}}\ 1]^{\mathsf{T}}$, the problem (9) can be transformed into homogeneous separable QCQP formulation as follows

$$\min_{\mathbf{z}} \quad \mathbf{z}^{\mathsf{T}} \mathbf{B}\, \mathbf{z} \tag{10a}$$

$$\text{s.t.} \quad \mathbf{z}^{\mathsf{T}} \mathbf{B}_{x,m}\, \mathbf{z} = 0, \tag{10b}$$

$$\mathbf{z}^{\mathsf{T}} \mathbf{B}_{y,m}\, \mathbf{z} = 0, \tag{10c}$$

$$\mathbf{z}^{\mathsf{T}} \mathbf{B}_{xy,m}\, \mathbf{z} = 1, \tag{10d}$$

$$\sum_{m=1}^{M} \mathbf{z}^{\mathsf{T}} \mathbf{B}_{\beta,m}\, \mathbf{z} \leq 1, \tag{10e}$$

$$\mathbf{z}^{\mathsf{T}} \mathbf{B}_{\beta xy,m}\, \mathbf{z} \leq 0, \tag{10f}$$

where

$$\mathbf{B} = \begin{bmatrix} \mathbf{0}_{(3M+1)\times(3M+1)} & \frac{1}{2}\mathbf{e}_{(3M+1)} \\ \frac{1}{2}\mathbf{e}_{(3M+1)}^{\mathsf{T}} & 0 \end{bmatrix},$$

$$\mathbf{B}_{y,m} = \begin{bmatrix} \mathbf{A}_{y,m} & -\frac{1}{2}\mathbf{e}_{M+m} \\ -\frac{1}{2}\mathbf{e}_{M+m}^{\mathsf{T}} & 0 \end{bmatrix}, \ \mathbf{b}_{xy,m} = \mathbf{e}_{m} + \mathbf{e}_{m+M},$$

$$\mathbf{B}_{xy,m} = \begin{bmatrix} \mathbf{0}_{(3M+1)\times(3M+1)} & \frac{1}{2}\mathbf{b}_{xy,m} \\ \frac{1}{2}\mathbf{b}_{xy,m}^{\mathsf{T}} & 0 \end{bmatrix},$$

$$\mathbf{B}_{\beta,m} = \begin{bmatrix} \mathbf{0}_{(3M+1)\times(3M+1)} & \frac{1}{2}\mathbf{e}_{m+2M} \\ \frac{1}{2}\mathbf{e}_{m+2M}^{\mathsf{T}} & 0 \end{bmatrix},$$

$$\mathbf{B}_{\beta xy,m} = \begin{bmatrix} \mathbf{A}_{\beta xy,m} & \frac{1}{2}\mathbf{b}_{cy,m} \\ \frac{1}{2}\mathbf{b}_{cy,m}^{\mathsf{T}} & 0 \end{bmatrix}, \mathbf{B}_{x,m} = \begin{bmatrix} \mathbf{A}_{x,m} & \frac{1}{2}\mathbf{e}_{m} \\ \frac{1}{2}\mathbf{e}_{m}^{\mathsf{T}} & 0 \end{bmatrix},$$

$$\mathbf{A}_{\beta xy,m} = \mathbf{A}_{\beta x,m} + \mathbf{A}_{\beta y,m} + \mathbf{A}_{\beta t,m}, \ \mathbf{b}_{cy,m} = k_{m}^{c}\mathbf{e}_{m+M}.$$

Next, we apply the SDR to obtain the desired results. Let $\mathbf{Y} = \mathbf{z}\,\mathbf{z}^{\mathsf{T}}$ with $\mathrm{rank}(\mathbf{Y}) = 1$. Then, the separable semi-definite programming (SDP) problem can be expressed by relaxing problem (10) is as follows

$$\min_{\mathbf{Y}} \quad \mathrm{Tr}(\mathbf{B}\,\mathbf{Y}) \tag{11a}$$

$$\text{s.t.} \quad \mathrm{Tr}(\mathbf{B}_{x,m}\,\mathbf{Y}) = 0, \tag{11b}$$

$$\mathrm{Tr}(\mathbf{B}_{y,m}\,\mathbf{Y}) = 0, \tag{11c}$$

$$\mathrm{Tr}(\mathbf{B}_{xy,m}\,\mathbf{Y}) = 1, \tag{11d}$$

$$\sum_{m=1}^{M} \mathrm{Tr}(\mathbf{B}_{\beta,m}\mathbf{Y}) \leq 1, \tag{11e}$$

$$\mathrm{Tr}(\mathbf{B}_{\beta xy,m}\mathbf{Y}) \leq 0. \tag{11f}$$

We solve the above SDP problem in a polynomial time using a standard SDP software SeDuMi [19]. We obtain the offloading decision $x_m$ and $y_m$ of the original problem (8) from $\mathbf{Y}$. We use randomization method [20] to find binary offloading decisions. Accordingly, the probability of task processing at end-user and edge server is given as

$$P_m^l = \frac{p_m^l}{p_m^l(1 - p_m^e) + (1 - p_m^l)\,p_m^e}, \tag{12a}$$

$$P_m^e = \frac{p_m^e}{p_m^l(1 - p_m^e) + (1 - p_m^l)\,p_m^e}, \tag{12b}$$

where $p_m^l = x_m$ and $p_m^e = y_m$. Now, we generate $N$ i.i.d. feasible offloading solutions as $\boldsymbol{\xi}^{(n)} = [(q_1^{(n)})^{\mathsf{T}} \ldots (q_M^{(n)})^{\mathsf{T}}]^{\mathsf{T}}$ using the probabilities in (12), for $n = 1, \ldots, N$, as follows

$$q_m = \begin{cases} [1,\, 0] & \text{with probability } P_m^l \text{ (at local CPU)}, \\ [0,\, 1] & \text{with probability } P_m^e \text{ (at edge server)}. \end{cases} \tag{13}$$

Next, we solve the problem (5) for the optimal resource allocation corresponding to offloading decision $\boldsymbol{\xi}^{(n)}$ obtained using (13). Therefore, (4) can be rewritten as

$$T_i^{\mathsf{L}} + T_i^{\mathsf{E}} = k_i^f + \frac{k_i^\eta}{\beta_i}, \tag{14a}$$

$$T_j^{\mathsf{L}} + T_j^{\mathsf{E}} = k_j^f + \frac{k_j^\eta}{\beta_j}, \tag{14b}$$

where $k_i^f = (D_i\,L\,x_i)/f_i^l + (D_i\,L\,y_i)/f_i^e$, $k_i^\eta = (D_i\,y_i)/(\eta_i\,C)$, $k_j^f = (D_j\,L\,x_j)/f_j^l + (D_j\,L\,y_j)/f_j^e$, and $k_j^\eta = (D_j\,y_j)/(\eta_j\,C)$. Hence, our optimization problem becomes

$$\min_{\mathbf{r}} \quad \max\{(T_i^{\mathsf{L}} + T_i^{\mathsf{E}}), (T_j^{\mathsf{L}} + T_j^{\mathsf{E}})\} \ \forall\ i \in \mathcal{N}_1, j \in \mathcal{N}_2 \tag{15a}$$

$$\text{s.t.} \quad \sum_{i=1}^{K} y_i\beta_i + \sum_{j=K+1}^{M} y_j\beta_j \leq 1. \tag{15b}$$

Now, we take an auxiliary variable $\theta$ as

$$\max_{i \in \mathcal{N}_1, j \in \mathcal{N}_2} \left\{(T_i^{\mathsf{L}} + T_i^{\mathsf{E}}), (T_j^{\mathsf{L}} + T_j^{\mathsf{E}})\right\} = \theta, \tag{16}$$

and from (14) and (16), we can write

$$k_i^f\,\beta_i + k_i^\eta - \beta_i\,\theta \leq 0, \tag{17a}$$

$$k_j^f\,\beta_j + k_j^\eta - \beta_j\,\theta \leq 0. \tag{17b}$$

Then, the optimization problem becomes

$$\min_{\mathbf{r}} \quad \theta \tag{18a}$$

$$\text{s.t.} \quad \text{(15b), (17a), and (17b)}. \tag{18b}$$

### C. Vector-matrix Formation

We denote $\mathbf{v} = [\beta_1, \beta_2, \ldots, \beta_K, \beta_{K+1}, \ldots, \beta_M, \theta, 1]^{\mathsf{T}}$. Defining a unit vector as $\hat{\mathbf{u}}_p = [\mathbf{0}_{1\times(p-1)}, 1, \mathbf{0}_{1\times(M+2-p)}]^{\mathsf{T}}$, the vector-matrix form of problem (18) is written as

$$\min_{\mathbf{v}} \quad \hat{\mathbf{u}}_{(M+1)}^{\mathsf{T}}\mathbf{v} \tag{19a}$$

$$\text{s.t.} \quad \sum_{i=1}^{K} \mathbf{b}_{yu,i}^{\mathsf{T}}\mathbf{v} + \sum_{j=K+1}^{M} \mathbf{b}_{yu,j}^{\mathsf{T}}\mathbf{v} \leq 1, \tag{19b}$$

$$\mathbf{b}_{kf,i}^{\mathsf{T}}\mathbf{v} + \mathbf{b}_{k\eta,i}^{\mathsf{T}}\mathbf{v} + \mathbf{v}^{\mathsf{T}}A_{\beta\theta,i}\mathbf{v} \leq 0, \tag{19c}$$

$$\mathbf{b}_{kf,j}^{\mathsf{T}}\mathbf{v} + \mathbf{b}_{k\eta,j}^{\mathsf{T}}\mathbf{v} + \mathbf{v}^{\mathsf{T}}A_{\beta\theta,j}\mathbf{v} \leq 0, \tag{19d}$$

where

$$\mathbf{b}_{kf,i} = k_i^f \hat{\mathbf{u}}_i, \ \mathbf{b}_{k\eta,i} = k_i^\eta \hat{\mathbf{u}}_{M+2}, \ \mathbf{b}_{k,i} = \mathbf{b}_{kf,i} + \mathbf{b}_{k\eta,i},$$
$$\mathbf{b}_{yu,i} = y_i \hat{\mathbf{u}}_i, \ \mathbf{b}_{kf,j} = k_j^f \hat{\mathbf{u}}_j, \ \mathbf{b}_{k\eta,j} = k_j^\eta \hat{\mathbf{u}}_{M+2},$$
$$\mathbf{b}_{k,j} = \mathbf{b}_{kf,j} + \mathbf{b}_{k\eta,j}, \ \mathbf{b}_{yu,j} = y_j \hat{\mathbf{u}}_j ,$$

$$\mathbf{A}_{\beta\theta,i} = -\frac{1}{2} \begin{bmatrix} \mathbf{0}_{(i-1)\times(M+2)} \\ \hat{\mathbf{u}}_{M+1}^\mathsf{T} \\ \mathbf{0}_{(M-i)\times(M+2)} \\ \hat{\mathbf{u}}_i^\mathsf{T} \\ \mathbf{0}_{(1)\times(M+2)} \end{bmatrix},$$

$$\mathbf{A}_{\beta\theta,j} = -\frac{1}{2} \begin{bmatrix} \mathbf{0}_{(j-1)\times(M+2)} \\ \hat{\mathbf{u}}_{M+1}^\mathsf{T} \\ \mathbf{0}_{(M-j)\times(M+2)} \\ \hat{\mathbf{u}}_j^\mathsf{T} \\ \mathbf{0}_{(1)\times(M+2)} \end{bmatrix}.$$

Let $\mathbf{s}^\mathsf{T} = [\mathbf{v}^\mathsf{T} \ 1]^\mathsf{T}$, thus the objective function becomes

$$\min_{\mathbf{s}} \quad \mathbf{s}^\mathsf{T} \mathbf{H} \mathbf{s} \tag{20a}$$

$$\text{s.t.} \quad \sum_{i=1}^{K} \mathbf{s}^\mathsf{T} \mathbf{H}_{\beta,i} \mathbf{s} + \sum_{j=K+1}^{M} \mathbf{s}^\mathsf{T} \mathbf{H}_{\beta,j} \mathbf{s} \le 1, \tag{20b}$$

$$\mathbf{s}^\mathsf{T} \mathbf{H}_{\beta k\theta,i} \mathbf{s} \le 0, \tag{20c}$$

$$\mathbf{s}^\mathsf{T} \mathbf{H}_{\beta k\theta,j} \mathbf{s} \le 0, \tag{20d}$$

where

$$\mathbf{H} = \begin{bmatrix} \mathbf{0}_{(M+2)\times(M+2)} & \frac{1}{2}\hat{\mathbf{u}}_{(M+1)} \\ \frac{1}{2}\hat{\mathbf{u}}_{(M+1)}^\mathsf{T} & 0 \end{bmatrix},$$

$$\mathbf{H}_{\beta,i} = \begin{bmatrix} \mathbf{0}_{(M+2)\times(M+2)} & \frac{1}{2}\mathbf{b}_{yu,i} \\ \frac{1}{2}\mathbf{b}_{yu,i}^\mathsf{T} & 0 \end{bmatrix},$$

$$\mathbf{H}_{\beta,j} = \begin{bmatrix} \mathbf{0}_{(M+2)\times(M+2)} & \frac{1}{2}\mathbf{b}_{yu,j} \\ \frac{1}{2}\mathbf{b}_{yu,j}^\mathsf{T} & 0 \end{bmatrix},$$

$$\mathbf{H}_{\beta k\theta,i} = \begin{bmatrix} \mathbf{A}_{\beta\theta,i} & \frac{1}{2}\mathbf{b}_{k,i} \\ \frac{1}{2}\mathbf{b}_{k,i}^\mathsf{T} & 0 \end{bmatrix}, \ \mathbf{H}_{\beta k\theta,j} = \begin{bmatrix} \mathbf{A}_{\beta\theta,j} & \frac{1}{2}\mathbf{b}_{k,j} \\ \frac{1}{2}\mathbf{b}_{k,j}^\mathsf{T} & 0 \end{bmatrix}.$$

Further, applying the SDR to obtain the desired results, let $\mathbf{S} = \mathbf{s}\,\mathbf{s}^\mathsf{T}$ such that $\text{rank}(\mathbf{S}) = 1$. Then, the SDP problem, by relaxing problem (20), can be expressed as

$$\min_{\mathbf{S}} \quad \text{Tr}(\mathbf{H} \mathbf{S}) \tag{21a}$$

$$\text{s.t.} \quad \sum_{i=1}^{K} \text{Tr}(\mathbf{H}_{\beta,i} \mathbf{S}) + \sum_{j=K+1}^{M} \text{Tr}(\mathbf{H}_{\beta,j} \mathbf{S}) \le 1, \tag{21b}$$

$$\text{Tr}(\mathbf{H}_{\beta k\theta,i} \mathbf{S}) \le 0, \tag{21c}$$

$$\text{Tr}(\mathbf{H}_{\beta k\theta,j} \mathbf{S}) \le 0. \tag{21d}$$

We solve the above SDP problem (21) in a polynomial time, denoting $\mathbf{S}$ as the optimal solution of the SDP problem (21). Finally, we obtain the optimal values of $\beta_i$ and $\beta_j$ from $\mathbf{S}$.
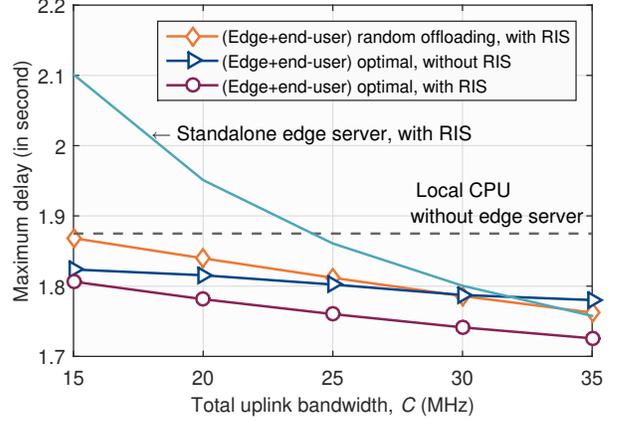


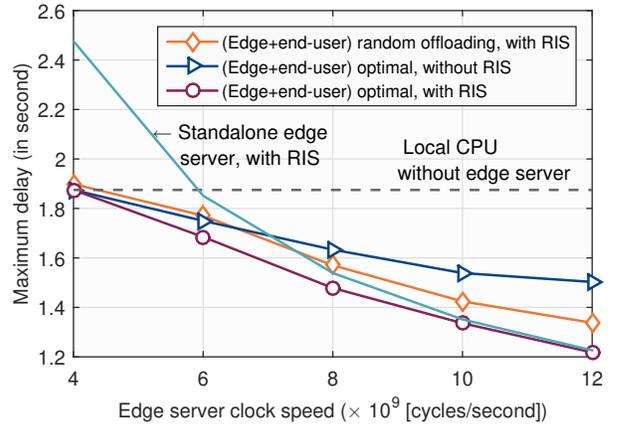Fig. 2. Delay performance vs. total uplink bandwidth.



Fig. 3. Delay performance with computing resource in edge server.

## IV. SIMULATION RESULTS

In this section, we evaluate the performance of the proposed RIS-assisted computation offloading policy with Monte Carlo simulations. Unless specified, we set CPU clock speed of end-user, i.e, $f_i^l = f_j^l = 500\times10^6$ [cycles/second] and edge server, i.e., $f^\mathsf{E} = 5\times10^9$ [cycles/second]. We assume that edge server equally distributes its CPU clock speed, $f^\mathsf{E}$, to every end-users, i.e., $f_i^e = f_j^e = f^\mathsf{E}/M$. We further consider that the size of a single task is uniformly distributed over $[0.1, 0.9]$ [MB] with a task processing density, $L = 1900$ [cycles/byte]. Also, the total available uplink bandwidth is $C = 15$ [MHz]. We assume that $M = 8$ end-users, out of which $(M - K) = 3$ end-users have poor wireless connection. Moreover, we set $\eta_i = 3.5$ [bps/Hz], $\eta_j = 0.1$ [bps/Hz] without RIS and $\eta_j = 3$ [bps/Hz] with RIS. We further set $N = 10$ as in [20]. The simulation results[1] are averaged over at least 10,000 different runs.

---

[1]The source code is available at https://github.com/MithunHub/GC2021Offloading.

Fig. 2 illustrates the maximum delay performance over the network with different uplink bandwidth. When we consider *'standalone edge server'*, the entire data is offloaded to the edge server. Thus, at lower uplink bandwidth, the worst performance is observed due to high uploading delay. From the figure, we can see that the maximum delay decreases with the increase of uplink bandwidth. The main reason is the uplink transmission delay decreases with the increase of uplink bandwidth. Therefore, the maximum delay over the network decreases. It is interesting to see that with RIS assistance, the maximum value of delay further reduces. Note that in this paper, we aim to minimize the maximum delay experienced by any end-user in the networks. Therefore, when no RIS support is available to the end-users with poor connection, the delay of these end-users has the adverse effect in minimizing the maximum value of delay in the network. Hence, reducing the uplink transmission delay with RIS assistance results in decreasing the maximum delay over the network.

Moreover, to show the performance with the computation resources in the edge server, Fig. 3 presents the maximum delay with increasing the CPU cycles/second of the edge server. From the figure, we observe that increase in CPU rate $f^E$ at the edge reduces the maximum delay. Moreover, one can clearly see that the offloading approach with RIS assistance exhibits better performance than the case without RIS assistance. In addition, we observe that at very high CPU cycles/second, the performance of standalone edge server and optimal offloading with RIS assistance gets very close to each others. Because, edge server's CPU rate is so high that end-users always prefer to offload the tasks under the setting of uplink data rate.

## V. Conclusions

In this paper, we studied computation offloading in a reconfigurable intelligent surface-assisted edge computing system. We employed the benefits of RIS to improve the uploading transmission rate for end-users with poor connection. Our proposed offloading scheme optimized the binary offloading decision variable, the uploading bandwidth allocation, and the CPU frequency allocated for the task data by the edge server. We applied SDR to solve the above QCQP problem. We note that with a better uplink quality, the poor user enjoys more chances to use the computational resources of the edge server, thereby improving the overall network performance than without RIS assistance. Our future work includes studying reliability and deadline constraints in task data offloading for an RIS-assisted edge computing system.

## Acknowledgment

## References

[1] T. Paikeday, "AI as you like it: NVIDIA DGX-ready partners make AI adoption easy," 2020, accessed on: July 20 2020. [Online]. Available: https://blogs.nvidia.com/blog/2020/05/14/dgx-ready-software-program-partners/

[2] C.-F. Liu, M. Bennis, M. Debbah, and H. V. Poor, "Dynamic task offloading and resource allocation for ultra-reliable low-latency edge computing," *IEEE Trans. Commun.*, vol. 67, no. 6, pp. 4132–4150, June 2019.

[3] J. Liu and Q. Zhang, "To improve service reliability for AI-powered time-critical services using imperfect transmission in MEC: An experimental study," *IEEE Internet of Things J.*, vol. 7, no. 10, pp. 9357–9371, Oct. 2020.

[4] M. Mukherjee, M. Guo, J. Lloret, and Q. Zhang, "Leveraging intelligent computation offloading with fog/edge computing for Tactile internet: Advantages and limitations," *IEEE Netw.*, vol. 34, no. 5, pp. 322–329, 2020.

[5] D. E. Boubiche, A.-S. K. Pathan, J. Lloret, H. Zhou, S. Hong, S. O. Amin, and M. A. Feki, "Advanced industrial wireless sensor networks and intelligent IoT," *IEEE Commun. Mag.*, vol. 56, no. 2, pp. 14–15, Feb. 2018.

[6] "Making boundless XR a commercial reality: Kicking off a trial to utilize existing 5G release-15 features to make XR available at scale," 2020, accessed on: July 20 2020. [Online]. Available: https://www.qualcomm.com/news/onq/2020/05/27/making-boundless-xr-commercial-reality

[7] "Edgeconnecx," https://www.edgeconnex.com, accessed on: July 20 2020.

[8] E. Basar, M. Di Renzo, J. De Rosny, M. Debbah, M. S. Alouini, and R. Zhang, "Wireless communications through reconfigurable intelligent surfaces," *IEEE Access*, vol. 7, pp. 116 753–116 773, 2019.

[9] T. Bai, C. Pan, Y. Deng, M. Elkashlan, A. Nallanathan, and L. Hanzo, "Latency minimization for intelligent reflecting surface aided mobile edge computing," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 11, pp. 2666–2682, Nov. 2020.

[10] T. Bai, C. Pan, H. Ren, Y. Deng, M. Elkashlan, and A. Nallanathan, "Resource allocation for intelligent reflecting surface aided wireless powered mobile edge computing in OFDM systems," *IEEE Transactions on Wireless Communications*, 2021, to be published.

[11] Z. Chu, P. Xiao, M. Shojafar, D. Mi, J. Mao, and W. Hao, "Intelligent reflecting surface assisted mobile edge computing for internet of things," *IEEE Wireless Commun. Lett.*, vol. 10, no. 3, pp. 619–623, Mar. 2021.

[12] Y. Liu, J. Zhao, Z. Xiong, D. Niyato, C. Yuen, C. Pan, and B. Huang, "Intelligent reflecting surface meets mobile edge computing: Enhancing wireless communications for computation offloading," *arXiv:2001.07449*, 2020.

[13] F. Zhou, C. You, and R. Zhang, "Delay-optimal scheduling for IRS-aided mobile edge computing," *IEEE Wireless Commun. Lett.*, vol. 10, no. 4, pp. 740–744, Apr. 2021.

[14] Y. Cao, T. Lv, Z. Lin, and W. Ni, "Delay-constrained joint power control, user detection and passive beamforming in intelligent reflecting surface-assisted uplink mmwave system," *IEEE Transactions on Cognitive Communications and Networking*, 2021.

[15] W. Liu, X. Zang, Y. Li, and B. Vucetic, "Over-the-air computation systems: Optimization, analysis and scaling laws," *IEEE Trans. Wireless Commun.*, vol. 19, no. 8, pp. 5488–5502, Aug. 2020.

[16] W. Fang, M. Fu, K. Wang, Y. Shi, and Y. Zhou, "Stochastic beamforming for reconfigurable intelligent surface aided over-the-air computation," in *Proc. IEEE GLOBECOM*, Dec. 2020, pp. 1–6.

[17] Z. Wang, Y. Shi, Y. Zhou, H. Zhou, and N. Zhang, "Wireless-powered over-the-air computation in intelligent reflecting surface-aided IoT networks," *IEEE Internet of Things J.*, vol. 8, no. 3, pp. 1585–1598, Feb. 2021.

[18] M. Mukherjee, V. Kumar, M. Guo, D. B. da Costa, Z. D. Ertugrul Basar, and W. K. Wong, "The interplay of reconfigurable intelligent surfaces and mobile edge computing in future wireless networks: A win-win strategy to 6G," *arXiv:2106.11784*, 2021. [Online]. Available: https://arxiv.org/abs/2106.11784

[19] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.1," http://cvxr.com/cvx, Mar. 2014.

[20] M.-H. Chen, M. Dong, and B. Liang, "Resource sharing of a computing access point for multi-user mobile cloud offloading with delay constraints," *IEEE Trans. on Mobile Comput.*, pp. 1–13, 2018.