



HAL
open science

Multiple resource allocation in multi-tenant edge computing via sub-modular optimization

Ayoub Ben-Ameur, Andrea Araldo, Tijani Chahed

► **To cite this version:**

Ayoub Ben-Ameur, Andrea Araldo, Tijani Chahed. Multiple resource allocation in multi-tenant edge computing via sub-modular optimization. IEEE International Conference on Communications (ICC), IEEE, May 2023, Rome, Italy. pp.3738-3743, 10.48550/arXiv.2302.09888 . hal-04009368

HAL Id: hal-04009368

<https://hal.science/hal-04009368v1>

Submitted on 1 Mar 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Multiple Resource Allocation in Multi-Tenant Edge Computing via Sub-modular Optimization

Ayoub Ben-Ameur, Andrea Araldo, Tijani Chahed

SAMOVAR, Telecom SudParis, Institut Polytechnique de Paris, 91120 Palaiseau, France

{first_name}.{last_name}@telecom-sudparis.eu

Abstract—Edge Computing (EC) allows users to access computing resources at the network frontier, which paves the way for deploying delay-sensitive applications such as Mobile Augmented Reality (MAR). Under the EC paradigm, MAR users connect to the EC server, open sessions and send continuously frames to be processed. The EC server sends back virtual information to enhance the human perception of the world by merging it with the real environment. Resource allocation arises as a critical challenge when several MAR Service Providers (SPs) compete for limited resources at the edge of the network. In this paper, we consider EC in a multi-tenant environment where the resource owner, i.e., the Network Operator (NO), virtualizes the resources and lets SPs run their services using the allocated slice of resources. Indeed, for MAR applications, we focus on two specific resources: CPU and RAM, deployed in some edge node, e.g., a central office. We study the decision of the NO about how to partition these resources among several SPs. We model the arrival and service dynamics of users belonging to different SPs using Erlang queuing model and show that under perfect information, the interaction between the NO and SPs can be formulated as a sub-modular maximization problem under multiple Knapsack constraints. To solve the problem, we use an approximation algorithm, guaranteeing a bounded gap with respect to the optimal theoretical solution. Our numerical results show that the proposed algorithm outperforms baseline proportional allocation in terms of the number of sessions accommodated at the edge for each SP.

Index Terms—Resource allocation, multi-tenant edge computing, mobile augmented reality, multi-dimensional knapsack problem, queuing model.

I. INTRODUCTION

Mobile Augmented Reality (MAR) has become one of the most emerging applications, accompanied by the development of mobile devices and wireless communication. In MAR, the human perception of the world can be enhanced by merging virtual information (generated from object detection, classification, or tracking) with the real environment via mobile devices [1]. However, it is difficult for a mobile device to offer the abundant computation and energy required by MAR applications.

While the production of AR/VR dedicated hardware seems very effective to run AR/VR applications properly, it is costly in the sense that only big players can afford producing their own devices. Hence, multi-tenant EC is particularly interesting for all the other players, as it is probably the only way for small or medium AR Service Providers (SPs) to run their

applications at the edge of the network. The development of EC and 5G has eliminated the obstacle to deploying the MAR service. In the concept of EC [2], computing and storage resources are deployed at the edge of the access network. Several MAR clients on mobile devices can send MAR requests that contain original data captured by sensors and cameras to the Edge Computing (EC) server. Furthermore, dedicated computing hardware (e.g., Graphics Processing Unit (GPU) and Central Processing Unit (CPU)) and software (e.g., computer vision-based algorithms) process these data and then return the results, such as object classification or space coordinate information, to the mobile devices.

The use of the EC for MAR has attracted extensive attention from the research community and industry recently [3], [4], which mainly focus on architecture design and deployment. However, scheduling the MAR requests received from several competing MAR clients on one EC server is critical and challenging. We address in this work the issue of resource allocation to competing, heterogeneous SPs in the case of multiple, limited resources at the Edge. We first model the arrivals and service dynamics of the flows using Erlang queuing model. We then formulate the resource allocation problem to each of the SPs using sub-modular maximization under Knapsack constraints. We next propose an implementation of the so-called streaming algorithm to solve the allocation problem, and obtain a $(\frac{1}{1+2d} - \epsilon)$ -approximate optimal value, where d is the number of resource types and ϵ is a controllable error term. We eventually provide numerical results to show that the resulting system performance significantly outperforms baseline resource allocation policies.

The remainder of this paper is organized as follows. In Section II we discuss most relevant work related to ours. We introduce in Section III our system model. We formulate the sub-modular maximization problem under Knapsack constraints in Section IV and describe the proposed algorithm to solve it. In Section V, we show our simulation results. We draw conclusions in Section VI.

II. RELATED WORK

Recently, much research effort has been made to develop MAR applications under the EC paradigm. In addition to studies on efficient EC architecture design for MAR [5], [6], in preliminary studies, researchers concentrated on the

resource allocation problem in the MAR service [7], [8]. Some researchers began to notice the trade-off between processing latency and accuracy. They aimed to develop acceleration mechanisms to reduce processing latency [9] or characterize the relation between computational complexity and the image size [10]. Based on these studies, the adaptation of the client configuration (image size and frame rate), and the resource allocation scheme were jointly considered in a centralized manner [11], [12]. However, in both studies, the researchers ignored the characteristics of dedicated computing devices (i.e., using batch processing to improve the GPU utility) for MAR tasks. Moreover, their solutions centrally controlled each client configuration, which is challenging to apply to a realistic MAR system.

In [13], the authors consider an edge computing system under network slicing in which the wireless devices generate latency sensitive computational tasks. The allocation of wireless and computing resources to a set of autonomous wireless devices in an edge computing system is considered in [14]. A main common assumption of the papers above is that user devices submit tasks to the NO. Contention in these works is modeled among user devices. However, we consider that these models are not appropriate for EC in our vision, since all traffic between devices and service providers is encrypted to maintain confidentiality and the NO does not have control over it. Therefore the contention for resources is, in our vision, between SPs and not between tasks submitted by users. In our assumption, the NO can only decide how to allocate resources among SPs and then users device interact directly with SPs, outside the control of the NO. In [15], authors consider the interplay between latency constrained applications and function-level resource management in EC. A game theoretic model of the interaction between rate adaptive applications and a load balancing operator is developed under a function-oriented pay-as-you-go pricing model. In our approach, we assume that the NO does not require any payment from the SPs. The NO aims to maximize his own utility by allocating resources to SPs at the edge. In our vision, an important part of MAR providers cannot afford the payment for resources at the edge.

III. SYSTEM MODEL AND OPTIMIZATION PROBLEM

We consider a setting with one NO, owning a set of resources \mathcal{R} and willing to share them between P different SPs. Each SP can then use its assigned share as if it had a dedicated hardware deployed in the edge.

A. Request Pattern

MAR users of SP p arrive to the EC server following a Poisson process with rate λ_p expressed in *users/s*. Once a user of any SP p is connected to the edge server, a session is created. This session is valid for a period of time denoted T_p during which the user can perform a sequence of interactions within that MAR application. A single session can contain multiple activities all of which are stored in the session temporarily

while the user is connected. Each SP runs in the edge a virtual server, e.g., a Kubernetes POD [16]. A MAR user establishes a session with the virtual server of the respective SP. Within that session, it sends a stream of image processing requests. When users point their MAR device toward an object, raw video from the MAR device cameras are fetched and clipped into frames with specific image format, such as JPEG and PNG and sent to the edge server [5]. The video frames are delivered to the AR tracker to determine the user's position with respect to the physical surroundings. Given the tracking results, virtual coordinate of the environment can be established by the mapper. Then, the internal objects in video frames are identified by the object recognizer with robust features. The MAR device finally downloads information about the object from the edge server. The AR information is presented in a 3-D "experience" superimposed on the object. What users see, then, is part real and part virtual. Since MAR needs high data rates, ultra-low latency and the possible use of lightweight devices, performing processing at the edge of 5G mobile networks can help guarantee the requirements of MAR applications (Section III-F of [1]).

We assume that a session of a single user of SP p requires a certain amount of resource r denoted z_p^r . If the SP does not have at the edge such amount of resources available, the user will establish a session with the cloud, suffering longer delay. Once a user of SP p is served by the edge, his session will be closed and he leaves the EC system. Please note that users can leave the system when they decide, this does not deny that we can define an average service rate for SP p expressed in *users/s* denoted by $\mu_p = \frac{1}{T_p}$.

B. Resources Partitioning

The NO owns CPU and RAM at the edge of the network, for instance, in a server co-located with a (micro) base station or central offices at the metropolitan scale. It allocates a total capacity K^{CPU} of CPU and a total capacity K^{RAM} of RAM among the P SPs. The allocation is a vector $\vec{\theta} = (\theta^{\text{CPU}}, \theta^{\text{RAM}})$ where each vector θ^r is the allocation of resource r . More precisely, the allocation has a form as follows:

$$\vec{\theta} = (\theta_1^{\text{CPU}}, \dots, \theta_P^{\text{CPU}}, \theta_1^{\text{RAM}}, \dots, \theta_P^{\text{RAM}}) \quad (1)$$

We define the set of all possible allocations as:

$$\mathcal{T} \triangleq \left\{ \vec{\theta} \mid \sum_{p=1}^P \theta_p^r \leq K^r, \theta_p^r \in \mathbb{Z}^+, r \in \mathcal{R} \right\} \quad (2)$$

C. Service Model

We model our system as an Erlang queue [17] which models Poisson arrivals, exponentially distributed service time, and a number of servers equal to the number of places in the system, i.e., users are either directly served at the edge or directed to the cloud. In our case, users of SP p arrive to the edge according to a Poisson distribution with mean arrival rate λ_p , they remain in the system for an exponentially distributed duration, T_p . The

number of servers in our case refers to the maximum number of sessions that the edge can accommodate for each SP, as determined next. Each user of SP p has fixed requirements $(z_p^r)_{p=1..P, r \in \mathcal{R}}$ and fixed allocation $((\theta_p^r)_{r \in \mathcal{R}})$ during service. We denote by $n_p(\vec{\theta})$ the maximum number of users that can be served at the edge for a SP p when the resource allocation decided by the NO is $\vec{\theta} = (\theta_p^r)_{r \in \mathcal{R}}$. Each user of each SP p will receive an amount z_p^r of the resource r for their session. Hence the maximum number of sessions $n_p(\vec{\theta})$ each SP p can establish at the edge when the allocation from the NO is $\vec{\theta}$ must satisfy:

$$n_p(\vec{\theta}) \cdot z_p^r \leq \theta_p^r, p = 1 \dots P, r \in \mathcal{R}. \quad (3)$$

Therefore, $n_p(\vec{\theta})$ is:

$$n_p(\vec{\theta}) = \left\lfloor \min_{r \in \mathcal{R}} \left(\frac{\theta_p^r}{z_p^r} \right) \right\rfloor, p = 1 \dots P \quad (4)$$

where $\lfloor \cdot \rfloor$ is the floor function giving as output the greatest integer less than or equal to $\left(\frac{\theta_p^r}{z_p^r} \right)$.

Let us denote by N_p the number of users of SP p served at the edge if all the resources are allocated only to this SP p .

$$N_p = \left\lfloor \min_{r \in \mathcal{R}} \left(\frac{K^r}{z_p^r} \right) \right\rfloor, p = 1 \dots P \quad (5)$$

D. Utility Model

A user of SP p is served directly by the edge if the latter can satisfy the requirements z_p^{RAM} and z_p^{CPU} . Otherwise, the corresponding session is not accepted (we say that it is “blocked”, following the terminology from queuing theory) and directed to a remote cloud server. Using Erlang (equation (3.45) of [17]), the probability for a user of SP p to be blocked is

$$B_p(\vec{\theta}) = \frac{A_p^{n_p(\vec{\theta})}}{\sum_{i=0}^{n_p(\vec{\theta})} \frac{A_p^i}{i!}}, p = 1 \dots P \quad (6)$$

where $A_p = \frac{\lambda_p}{\mu_p}$. The probability for a user of SP p to have his/her session established with the edge is thus:

$$\bar{B}_p(\vec{\theta}) = 1 - B_p(\vec{\theta}). \quad (7)$$

The utility perceived by a user who establishes a session directly in the edge is U_E , while if the session is with the cloud, the utility is U_C . Such utilities take into account the impact on the Quality of Experience (QoE) of the delay to process every user request, accounting for a larger delay to reach the cloud. Hence, $U_E > U_C > 0$. For simplicity, we assume that U_E and U_C are the same for all SPs. Since $1 - B_p$ indicates the fraction of users of SP p establishing sessions with the edge,

the expected value of the utility perceived by a user of SP p is, by the theorem of total probability:

$$\begin{aligned} \mathbb{E}U_p(\vec{\theta}) &= \mathbb{P}(\text{session established with the edge}) \cdot U_E \\ &\quad + \mathbb{P}(\text{session established with the cloud}) \cdot U_C \\ &= \bar{B}_p(\vec{\theta}) \cdot U_E + (1 - \bar{B}_p(\vec{\theta})) \cdot U_C \\ &= (U_E - U_C) \cdot \bar{B}_p(\vec{\theta}) + U_C \end{aligned} \quad (8)$$

By the theorem of total expectation, the utility perceived by a generic user is

$$\begin{aligned} \mathbb{E}U(\vec{\theta}) &= \sum_{p=1}^P \mathbb{E}U_p(\vec{\theta}) \cdot \mathbb{P}(\text{new user is for SP } p) \\ &= \sum_{p=1}^P w_p \cdot \mathbb{E}U_p(\vec{\theta}) \end{aligned} \quad (9)$$

where $w_p = \frac{\lambda_p}{\sum_{p'} \lambda_{p'}}$.

E. Optimization Problem

The NO aims to maximize the expected value of the utility perceived by a generic user:

$$\begin{aligned} \max_{\vec{\theta}} \quad & \mathbb{E}U(\vec{\theta}) \\ \text{s.t.} \quad & \sum_{p=1}^P \theta_p^r \leq K^r, \forall r \in \mathcal{R} \end{aligned} \quad (10)$$

Replacing $\mathbb{E}U_p(\vec{\theta})$ with its value found in (8) and observing that $(U_E - U_C)$ and U_C are positive constants, the optimization problem becomes:

$$\begin{aligned} \max_{\vec{\theta}} \quad & \sum_{p=1}^P w_p \bar{B}_p(\vec{\theta}) \\ \text{s.t.} \quad & \sum_{p=1}^P \theta_p^r \leq K^r, \forall r \in \mathcal{R} \end{aligned} \quad (11)$$

Thanks to (3) and (6), we can express the problem in terms of $\vec{n} = (n_1, \dots, n_P)$ instead of $\vec{\theta}$:

$$\begin{aligned} \max_{\vec{n}} \quad & f(\vec{n}) = \sum_{p=1}^P w_p \bar{B}_p(\vec{n}) \\ \text{s.t.} \quad & \sum_{p=1}^P n_p \cdot z_p^r \leq K^r, \forall r \in \mathcal{R} \end{aligned} \quad (12)$$

$$\text{where } \bar{B}_p(\vec{n}) \triangleq 1 - \frac{A_p^{n_p}}{\sum_{i=0}^{n_p} \frac{A_p^i}{i!}}, p = 1 \dots P \quad (13)$$

Observe that $f(\vec{n})$ is the probability for a generic user to be served with a session with the edge node. This shows that improving the expected user utility (10) is equivalent to maximizing the probability of establishing a session with the edge (12).

IV. SUB-MODULAR OPTIMIZATION

To describe our problem (12) in terms of sub-modular optimization, we interpret a user session established with the edge node as an item. Let $\mathcal{V}_p = \{1, 2, \dots, N_p\}$ be the set of *candidate sessions* of SP p that could coexist in the edge if all resources were given to this SP p . Since in reality resources at the edge are not given to one SP only, we need to choose a subset of sessions $\mathcal{S}_p \subseteq \mathcal{V}_p$ to allocate to each SP p . This choice induces a certain probability of establishing a session with the edge:

$$\bar{B}_p(\mathcal{S}_p) = 1 - \frac{A_p^{|\mathcal{S}_p|}}{|\mathcal{S}_p|!} \quad (14)$$

With slight abuse of notation, in the formula above we use the notation $\bar{B}_p(\cdot)$ as in (13), to emphasize that the two quantities are conceptually the same thing, by setting $n_p = |\mathcal{S}_p|$. Let $\mathcal{V} \triangleq \bigcup_{p=1}^P \mathcal{V}_p$ the set of all candidate sessions and $\mathcal{S} = \bigcup_{p=1}^P \mathcal{S}_p \subseteq \mathcal{V}$ the set of sessions allocated. Set \mathcal{S} is our decision variable. For each SP p , we define a non-negative set function f_p , taking as input all possible subsets \mathcal{S} of \mathcal{V} , as follows:

$$f_p(\mathcal{S}) \triangleq w_p \cdot \bar{B}_p(\mathcal{S} \cap \mathcal{V}_p) \in [0, 1]$$

Function f_p represents the probability, for a user that arrives, to be of SP p and to be served with a session at the edge. We define $f(\mathcal{S}) \triangleq \sum_{p=1}^P f_p(\mathcal{S})$. It indicates, for any arriving user, the probability to be served with a session at the edge.

For any subset \mathcal{S} of \mathcal{V} , we denote the characteristic vector of \mathcal{S} by $\mathbf{x}_{\mathcal{S}} = (x_{\mathcal{S}_1,1}, \dots, x_{\mathcal{S}_1,N_1}, \dots, x_{\mathcal{S}_P,1}, \dots, x_{\mathcal{S}_P,N_P})^T$, where for any $j \in [1, N_p]$ and $p = 1, \dots, P$:

$$x_{\mathcal{S}_p,j} = \begin{cases} 1, & \text{if the } j\text{-th item of } \mathcal{V}_p \text{ is in } \mathcal{S}_p \\ 0, & \text{otherwise} \end{cases}$$

For $\mathcal{S} \subseteq \mathcal{V}$ and $v \in \mathcal{V}$, the marginal gain in f when adding v to set \mathcal{S} is defined as $\Delta_f(v|\mathcal{S}) \triangleq f(\mathcal{S} \cup \{v\}) - f(\mathcal{S})$.

We introduce now the d -knapsack constraint where $d = |\mathcal{R}|$. Let $\mathbf{k} = (K^1, \dots, K^d)^T$ be the resource capacity vector and $\mathbf{Z}_p = (z_{p,j}^r)$ denote a $d \times N_p$ matrix, whose (r, j) -th entry $z_{p,j}^r > 0$ is the weight of the j -th item of \mathcal{V}_p in terms of resource r . Since we have assumed (§III-C) that all users of a SP p require the same amount of each resource, $z_{p,j}^r = z_p^r$ for all the items in \mathcal{V}_p . Therefore, the constraint in (12) can be expressed by $\mathbf{Z} \cdot \mathbf{x}_{\mathcal{S}} \leq \mathbf{k}$, where $\mathbf{Z} = (\mathbf{Z}_1, \dots, \mathbf{Z}_P) \in \mathbb{R}^{d \times \sum_p N_p}$ and $\mathbf{x}_{\mathcal{S}} \in \{0, 1\}^{\sum_p N_p \times 1}$. Problem (12) becomes:

$$\begin{aligned} \max_{\mathcal{S}} \quad & f(\mathcal{S}) = \sum_{p=1}^P f_p(\mathcal{S}_p) \\ \text{s.t.} \quad & \mathbf{Z} \mathbf{x}_{\mathcal{S}} \leq \mathbf{k} \end{aligned} \quad (15)$$

Without loss of generality, for $1 \leq i \leq d, 1 \leq j \leq N$, we assume that $z_p^r \leq K^r$. That is, no item has a larger weight

than the corresponding knapsack budget, since otherwise such an item would never be selected into \mathcal{S} .

We are now ready to study the properties of formulation (15). To do so, we recall two common definitions from set-function theory [18].

Algorithm 1: Streaming Algorithm for sub-modular maximization problem under Knapsack constraints

Data: $d, z_p^r, K^r, \lambda_p, \mu_p$

Result: \mathcal{S}^*

$m \leftarrow 0;$

$\mathcal{Q} \leftarrow \{[1 + (1 + 2d)\epsilon]^l | l \in \mathbb{Z}\};$

for $v \in \mathcal{Q}$ **do**

$\mathcal{S}_v \leftarrow \emptyset;$

for $1 \leq i \leq d$ **do**

$m \leftarrow \max\{m, f(\{j\})/z_{i,j}\};$

end

$\mathcal{Q} \leftarrow \{[1 + (1 + 2d)\epsilon]^l | l \in \mathbb{Z},$

$\frac{m}{1+(1+2d)\epsilon} \leq [1 + (1 + 2d)\epsilon]^l \leq 2Km\};$

for $1 \leq j \leq n$ **do**

if $\exists i \in [1, d], z_{i,j} \geq \frac{K}{2}$ **and** $\frac{f(\{j\})}{z_{i,j}} \geq \frac{2v}{K^{(1+2d)}}$

then

$\mathcal{S}_v \leftarrow \{j\};$

break;

end

if $\forall i \in [1, d], \sum_{l \in \mathcal{S} \cup \{j\}} z_{i,l} \leq K$ **and** $\frac{\Delta_f(j|\mathcal{S})}{z_{i,j}} \geq$

$\frac{2v}{K^{(1+2d)}}$ **then**

$\mathcal{S}_v \leftarrow \mathcal{S}_v \cup \{j\};$

end

end

end

$\mathcal{S}^* \leftarrow \arg \max_{\mathcal{S}_v, v \in \mathcal{Q}} f(\mathcal{S}_v);$

Definition IV-1. A function f is sub-modular if it satisfies that $\Delta_f(v|\mathcal{B}) \leq \Delta_f(v|\mathcal{A})$, for any $\mathcal{A} \subseteq \mathcal{B} \subseteq \mathcal{V}$ and $v \in \mathcal{V} \setminus \mathcal{B}$.

Definition IV-2. A function f is monotone if for any $\mathcal{S} \subseteq \mathcal{V}$ and $v \in \mathcal{V}$, $\Delta_f(v|\mathcal{S}) \geq 0$.

Theorem IV-3. Function f in (15) is monotone and sub-modular.

Proof. Let $\mathcal{S} \subseteq \mathcal{V}$ and $v \in \mathcal{V}$. Suppose in particular that $v \in \mathcal{V}_{p'}$.

$$\begin{aligned} \Delta_f(v|\mathcal{S}) &= f(\mathcal{S} \cup \{v\}) - f(\mathcal{S}) \\ &= \sum_{p \neq p'} f_p(\mathcal{S}_p) + f_{p'}(\mathcal{S}_{p'} \cup \{v\}) - \sum_{p=1}^P f_p(\mathcal{S}_p) \\ &= f_{p'}(\mathcal{S}_{p'} \cup \{v\}) - f_{p'}(\mathcal{S}_{p'}) \\ &= w_{p'} \cdot \bar{B}_{p'}(\mathcal{S}_{p'} \cup \{v\}) - w_{p'} \cdot \bar{B}_{p'}(\mathcal{S}_{p'}) \\ &\geq 0, \end{aligned}$$

where the last inequality can be obtained by simple calculus from (14). This shows that function f is monotone.

Let us consider sets $\mathcal{A} \subseteq \mathcal{B} \subseteq \mathcal{V}$ and a vector $v \in \mathcal{V} \setminus \mathcal{B}$.

$$\begin{aligned} \Delta_f(v|\mathcal{B}) - \Delta_f(v|\mathcal{A}) &= [f(\mathcal{B} \cup \{v\}) - f(\mathcal{B})] - \\ &\quad [f(\mathcal{A} \cup \{v\}) - f(\mathcal{A})] \\ &= [f(\mathcal{B} \cup \{v\}) - f(\mathcal{A} \cup \{v\})] + \\ &\quad [f(\mathcal{A}) - f(\mathcal{B})] \end{aligned}$$

Having $\mathcal{A} \subseteq \mathcal{B}$, we can write $\exists \Omega \subseteq \mathcal{V}/\mathcal{B} = \mathcal{A} \cup \Omega$. Hence:

$$\begin{aligned} &[f(\mathcal{B} \cup \{v\}) - f(\mathcal{A} \cup \{v\})] + [f(\mathcal{A}) - f(\mathcal{B})] \\ &= [f(\mathcal{A} \cup \Omega \cup \{v\}) - f(\mathcal{A} \cup \{v\})] + [f(\mathcal{A}) - f(\mathcal{A} \cup \Omega)] \\ &\leq [f(\mathcal{A}) + f(\Omega \cup \{v\}) - f(\mathcal{A}) - f(\{v\})] \\ &\quad + [f(\mathcal{A}) - f(\mathcal{A} \cup \Omega)] \\ &\leq [f(\Omega \cup \{v\}) - f(\{v\})] + [f(\mathcal{A}) - f(\mathcal{A}) - f(\Omega)] \\ &= f(\Omega \cup \{v\}) - [f(\{v\}) + f(\Omega)] \leq 0 \end{aligned}$$

Therefore, the function f is sub-modular. \square

Now that we have proved that our objective function f is monotone and sub-modular, we can use well known results from sub-modular optimization. IN particular, we adopt the algorithms proposed in [19], which we report in Algorithm 1. The main idea of the algorithm is for every potential new user for each SP p , we compare the increase in f when we add this user to the set of users \mathcal{S} . We add the user providing the most increase in f . The algorithm guarantees the following sub-optimality gap (Theorem 1 of [19]).

Theorem IV-4. *Algorithm 1 outputs \mathcal{S} that satisfies $f(\mathcal{S}) \geq (\frac{1}{1+2d} - \epsilon)OPT$ and has $O(\frac{\log(K_{\max})}{\epsilon})$ computational complexity per element, d being the number of resources, $0 < \epsilon < \frac{1}{1+2d}$, $K_{\max} = \max_{1 \leq i \leq d} K^i$ and OPT the value of f obtained by the optimal solution.*

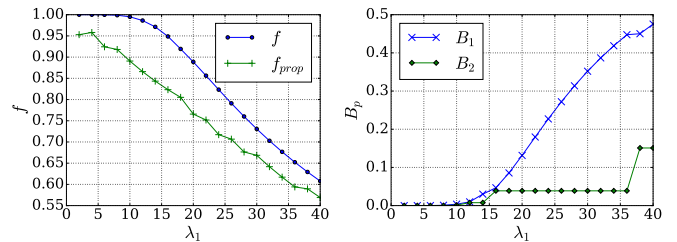
Note that the hyper-parameter ϵ impacts the behavior of the algorithm as well as the quality of the optimality gap. The smaller is ϵ , the larger is our $f(\mathcal{S})$.

V. NUMERICAL RESULTS

We now evaluate the performance of Algorithm 1 via a numerical model developed in Python and compare it to the proportional allocation where θ_p^r is proportional to the arrival rate λ_p of users of each SP p . We set $\epsilon = 0.01$.

A. Setting

We focus on an edge node co-located with a central offices serving 2 SPs. We set arrival rates λ_1 and λ_2 at 20 and 5 users/s, respectively and departure rates μ_1 and μ_2 at 1 and 10 users/s, respectively. Motivated by Amazon EC2 instances, such as G4dn [20], designed to support machine learning inference for applications like adding metadata to an image, object detection, recommendation systems, automated speech



(a) Objective function f vs. λ_1 (b) Blocking probability B_p vs. λ_1

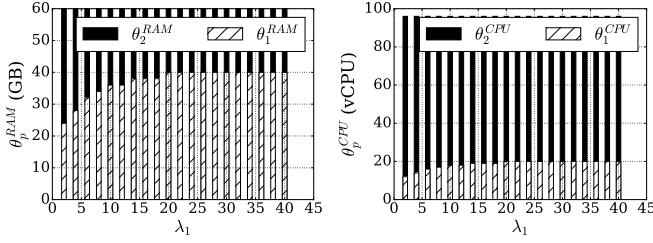
Fig. 1: Performance of the streaming algorithm w.r.t λ_1

recognition, and language translation, we consider an edge server similar to the G4dn.metal with $K^{\text{RAM}} = 384$ GB of total RAM capacity and a 2nd Generation Intel Xeon Scalable CPU: Cascade Lake P-8259L with total capacity of CPU $K^{\text{CPU}} = 96$ vCPU. Taking in consideration AR applications similar to Pokemon GO [21], we set RAM and CPU requirements for SP 1 and SP 2 at: $z_1^{\text{RAM}} = 2$ GB, $z_1^{\text{CPU}} = 1$ vCPU, $z_2^{\text{RAM}} = 0.5$ GB and $z_2^{\text{CPU}} = 4$ vCPU, respectively.

B. Results

We plot in Fig. 1a our solution obtained with Algorithm 1: the objective function f , which is the probability for a user to establish a session with the edge (12) and we compare our solution with the baseline f_{prop} , i.e., the probability of establishing sessions with the edge obtained when allocating resources to SPs proportionally to their users arrival rates. In Fig. 1b, we show the variation of the blocking probabilities for each SP when varying λ_1 . The increase in λ_1 results higher blocking probability for SP 1, which is expected as more users will consume more resources at the edge and less resources are left. Higher λ_1 will also affect SP 2 but much less significantly. As for resource utilization, we plot Fig. 2. The results show that the CPU is totally utilized by the two SPs (Fig. 2b), while the RAM is not fully exploited (less than 20% as shown in Fig. 2a). Despite having more than 80% of RAM free, we cannot expect better performance since the blocking comes always from the CPU, which is the scarcer resource. Having higher arrival rate, the algorithm does not allow yet SP 1 to have more CPU as this resource is almost 80% used by SP 2. We can explain this by looking to the values of z_1^{CPU} and z_2^{CPU} , we can see that SP 2 is CPU-greedy: users of SP 2 consume 4 times more CPU than users of SP 1.

In Fig.3, we plot a heat-map describing the global objective function f with respect to the variations of the two arrival rates. Obviously, the performance of the algorithm under lower arrival rates is better (dark red region $f \geq 0.95$). But what is more interesting in the figure, is that even for high arrival rates for SP 2 ($\lambda_2 \geq 35$), the algorithms keeps performing well up to $\lambda_1 = 20$ (orange region $f \geq 0.85$), no matter the arrival rate λ_2 of SP 2. The opposite is not the same: for any value of λ_2 , even small ones, the performance highly depend on λ_1 . We can



(a) RAM allocation vs. λ_1 (b) CPU allocation vs. λ_1

Fig. 2: Resource utilization vs. λ_1

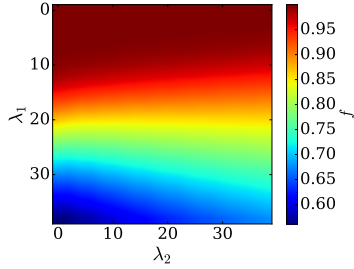


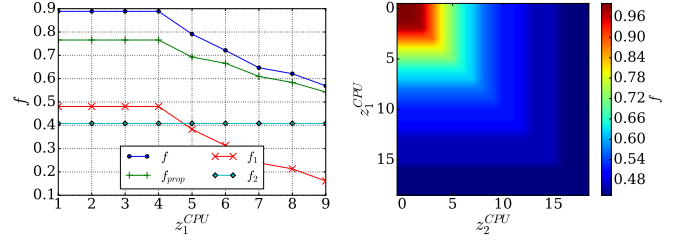
Fig. 3: Objective function f w.r.t λ_1 and λ_2

explain that by the fact that the users of SP 2 consume a lot of CPU (the blocking resource) which means every new admission of SP 1 would degrade the performance of the algorithm.

Since the CPU is the blocking resource, we evaluate in Fig. 4 the sensitivity of the system with respect to the required amount of CPU by each user of the two SPs. First, we plot in Fig. 4a the objective functions: f , f_1 and f_2 obtained by the algorithm and f_{prop} . The results show that the streaming algorithm outperforms the baseline allocation whatever users of SP 1 require in term of CPU. In Fig. 4b, we plot the heat-map describing the global objective function f obtained with the streaming algorithm with respect to the variations of the CPU requirements. The algorithm maintains a satisfying performance (dark red to light green region) up to requirements around 5 vCPU at most and then the performance rapidly decrease with the higher CPU requirements.

VI. CONCLUSION AND FUTURE WORK

We tackled in this paper resource allocation at EC between heterogeneous, MAR-oriented SPs competing over multiple, limited resources. We modeled the users dynamics in terms of an Erlang-type queuing model, we formulated the resource allocation problem as a sub-modular maximization problem subject to multiple knapsack constraints and solve it via an approximation algorithm with provable optimality gap. Our numerical results quantified the performance of our algorithm in terms of the probability that users get served by the Edge, as opposed to being blocked and re-directed towards the Cloud which entails larger delay and hence lesser QoE. We showed the resulting resources partitioning between the SPs.



(a) f vs. z_1^{CPU} (b) f w.r.t z_1^{CPU} and z_2^{CPU}

Fig. 4: Sensitivity w.r.t z_p^{CPU} , $p = 1, 2$

We showed the algorithm outperforms a baseline resource allocation, proportional to users arrival rates. Finally, we included a sensitivity analysis with respect to individual user requirement of a given resource. Our next work perspective would focus on the case where users arrival rates as well as resource requirements are unknown, the NO shall implement learning in order to be able to allocate resources in this case.

REFERENCES

- [1] Y. Siriwardhana *et al.*, "A survey on mobile augmented reality with 5G mobile edge computing: Architectures, applications, and technical aspects," *IEEE Communications Surveys & Tutorials*, 2021.
- [2] Y. Mao *et al.*, "A survey on mobile edge computing: The communication perspective," *IEEE communications surveys & tutorials*, 2017.
- [3] M. Erol-Kantarci *et al.*, "Caching and computing at the edge for mobile (AR/VR) in 5G," *Ad Hoc Networks*, 2018.
- [4] A. B. Ameer *et al.*, "On the deployability of augmented reality using embedded edge devices," in *IEEE CCNC*, 2021.
- [5] J. Ren *et al.*, "An edge-computing based architecture for mobile augmented reality," *IEEE Network*, 2019.
- [6] T. M. Fernández-Caramés *et al.*, "A fog computing and cloudlet based augmented reality system for the industry 4.0 shipyard," *Sensors*, 2018.
- [7] W. Liu *et al.*, "Data offloading and sharing for latency minimization in augmented reality based on mobile-edge computing," in *IEEE VTC*, 2018.
- [8] M. Jia and W. Liang, "Delay-sensitive multiplayer augmented reality game planning in mobile edge computing," in *ACM ICMAWSMS*, 2018.
- [9] N. Lane *et al.*, "DeepX: A software accelerator for low-power deep learning inference on mobile devices," in *ACM/IEEE IPSN*, 2016.
- [10] Y. He *et al.*, "Optimizing the learning performance in mobile augmented reality systems with cnn," *ToWC*, 2020.
- [11] Q. Liu *et al.*, "An edge network orchestrator for mobile augmented reality," in *IEEE INFOCOM 2018*, 2018.
- [12] Q. Liu and T. Han, "Dare: Dynamic adaptive mobile augmented reality with edge computing," in *IEEE ICNP*, 2018.
- [13] S. Jošilo *et al.*, "Joint wireless and edge computing resource management with dynamic network slice selection," *IEEE/ACM ToN*, 2022.
- [14] —, "Wireless and computing resource allocation for selfish computation offloading in edge computing," in *IEEE INFOCOM*, 2019.
- [15] F. Tütüncüoğlu *et al.*, "Online learning for rate-adaptive task offloading under latency constraints in serverless ec," *IEEE/ACM ToN*, 2022.
- [16] T. K. Authors. Kubernetes documentation. [Online]. Available: <https://kubernetes.io/docs/concepts/workloads/pods/>
- [17] L. Kleinrock, *Queuing Systems*. Wiley-Interscience, 1975, vol. 1.
- [18] S. Fujishige, *Submodular functions and optimization*. Elsevier, 2005.
- [19] Q. Yu *et al.*, "Submodular maximization with multi-knapsack constraints and its applications in scientific literature recommendations," in *IEEE GlobalSIP*, 2016.
- [20] (2022) Types of instances in amazon EC2. [Online]. Available: <https://aws.amazon.com/fr/ec2/instance-types/>
- [21] (2022) Pokemon GO requirements. [Online]. Available: <https://support.pokemon.com/hc/en-us/articles/-/Pokemon-GO-Plus-system-requirements-and-compatibility>