

Joint Wireless and Edge Computing Resource Management with Dynamic Network Slice Selection

Sladana Jošilo and György Dán

Division of Network and Systems Engineering,

School of Electrical Engineering and Computer Science

KTH, Royal Institute of Technology, Stockholm, Sweden E-mail: {josilo, gyuri}@kth.se

Abstract—Network slicing is a promising approach for enabling low latency computation offloading in edge computing systems. In this paper, we consider an edge computing system under network slicing in which the wireless devices generate latency sensitive computational tasks. We address the problem of joint dynamic assignment of computational tasks to slices, management of radio resources across slices and management of radio and computing resources within slices. We formulate the *Joint Slice Selection and Edge Resource Management (JSS-ERM)* problem as a mixed-integer problem with the objective to minimize the completion time of computational tasks. We show that the JSS-ERM problem is NP-hard and develop an approximation algorithm with bounded approximation ratio based on a game theoretic treatment of the problem. We provide extensive simulation results to show that network slicing can improve the system performance compared to no slicing and that the proposed solution can achieve significant gains compared to the equal slicing policy. Our results also show that the computational complexity of the proposed algorithm is approximately linear in the number of devices.

I. INTRODUCTION

Network slicing is emerging as an enabler for providing logical networks that are customized to meet the needs of different kinds of applications, mostly in 5G mobile networks. Horizontal network slices are designed for specific classes of applications, e.g., streaming visual analytics, real-time control, or media delivery, while vertical network slices are designed for specific industries. Slicing is expected to allow flexible and efficient end-to-end provisioning of bandwidth, composition of in-network processing, e.g., in the form of service chains composed of virtual network functions (VNF), and the allocation of dedicated computing resources. At the same time it provides performance isolation. Slicing is particularly appealing in combination with edge computing, as network slicing could allow low latency access to customized computing services located in edge clouds [1], [2].

Flexibility in network slicing is achieved through service orchestration. Orchestration focuses on the deployment and service-aware adaptation of VNFs and edge cloud services based on predicted workloads. Recent works in the area addressed the joint placement and routing of service function chains, formulated as a virtual network embedding problem [3], and the problem of joint resource

dimensioning and routing [4], [5]. Typical objectives are maximization of the service capacity or profit under physical (bandwidth and computational power) resource constraints, or the minimization of the energy consumption subject to satisfying service demand.

Common to the works on service orchestration is that they assume that each application is mapped to a specific slice deterministically, and assume a static resource pool per slice so as to ensure performance isolation [3], [4], [5]. A deterministic mapping is, however, not mandatory in practice. While there may be a designated (default) slice for every application, most proposed architectures for network slicing define a set of allowed slices, and the assignment of an application to a slice can be decided dynamically based on the current workload and SLA requirements [6]. The dynamic assignment of applications to slices thus results in a mixture of workloads in the slices, and consequently calls for flexibility in allocating resources to slices.

The importance of resource management across slices has been widely accepted in the case of the radio access network (RAN) [6]. Such inter-slice resource allocation should happen at short time scales, taking into account slice-level service level agreements (SLAs) and technological constraints (e.g., available RAN technology, such as 5G NR or WiFi-Lic). Recent work in the area has focused on system aspects of virtualizing RANs [7], and on the allocation of virtual resource block groups to slices so as to maximize efficiency [8], but has not considered of the potential impact of inter-slice resource management on service orchestration and on the dynamic assignment of applications to slices. It is thus so far unclear how to perform joint resource management within and across slices, considering the orchestration of communication and computing resources simultaneously.

In this paper we address the problem of joint dynamic slice selection, inter-slice radio resource management and intra-slice radio and computing resource management for latency sensitive workloads, and make three important contributions. First, we formulate the joint slice selection and edge resource management (JSS-ERM) problem, and show that it is NP-hard. Second, we analyze the optimal solution structure, and we develop an efficient approximation algorithm with bounded approximation

ratio inspired by a game theoretic treatment of the problem. Third, we provide extensive numerical results to show that the resulting system performance significantly outperforms baseline resource allocation policies.

The rest of the paper is organized as follows. Section II introduces the system model and Section III the problem formulation. Sections IV and V provide the analytical results, and Section VI shows numerical results. Section VII discusses related work and Section VIII concludes the paper.

II. SYSTEM MODEL

We consider a slicing enabled mobile backhaul including mobile edge computing (MEC) resources that serves a set $\mathcal{N}=\{1, 2, \dots, N\}$ of wireless devices (WDs) that generate computationally intensive tasks. WDs can offload their tasks through a set $\mathcal{A}=\{1, 2, \dots, A\}$ of access points (APs) to a set $\mathcal{C}=\{1, 2, \dots, C\}$ of edge clouds (ECs). APs and ECs form the set $\mathcal{E} \triangleq \mathcal{A} \cup \mathcal{C}$ of edge resources. We denote by $\mathcal{S}=\{1, 2, \dots, S\}$ the set of slices in the network, which include certain combinations of computing resources (e.g., CPUs, GPUs, NPUs and/or FPGAs), optimized for executing some types of tasks.

We characterize a task generated by WD i by the size D_i of the input data and by its complexity, which we define as the *expected* number of instructions required to perform the computation. Since the WDs and the slices may have different instruction set architectures, the number of instructions required to execute the same task may also differ. Hence, for a task generated by WD i we denote by L_i and $L_{i,s}$ the expected number of instructions required to perform the computation locally and in slice s , respectively. Similar to other works [9], [10], [11], we consider that D_i , L_i and $L_{i,s}$ can be estimated from measurements by applying the methods described in [12], [13], [14].

We consider that each WD i generates a computational task at a time; each task is atomic and can be either offloaded for computation or performed locally on the WD it was generated at. In the case of offloading, the WD will be assigned to exactly one slice $s \in \mathcal{S}$ and within the slice to exactly one AP $a \in \mathcal{A}$ and to exactly one EC $c \in \mathcal{C}$. Therefore, we define the set of feasible decisions for WD i as $\mathfrak{D}_i \triangleq \{i\} \cup \{(a, c, s) | a \in \mathcal{A}, c \in \mathcal{C}, s \in \mathcal{S}\}$ and we use variable $d_i \in \mathfrak{D}_i$ to indicate the decision for WD i 's task (i.e., $d_i = i$ indicates that WD i performs the task locally and $d_i = (a, c, s)$ indicates that WD i should offload its task through AP a to EC c in slice s). Furthermore, we define a decision vector $\mathbf{d} \triangleq (d_i)_{i \in \mathcal{N}}$ as the collection of the decisions of all WDs and we define the set $\mathfrak{D} \triangleq \times_{i \in \mathcal{N}} \mathfrak{D}_i$, i.e., the set of all possible decision vectors.

For a decision vector $\mathbf{d} \in \mathfrak{D}$ we define the set $O_{(a,s)}(\mathbf{d}) \triangleq \{i \in \mathcal{N} | d_i = (a, \cdot, s)\}$ of all WDs that use AP a in slice s and the set $O_a(\mathbf{d}) = \cup_{s \in \mathcal{S}} O_{(a,s)}(\mathbf{d})$ of all WDs that use AP a . Similarly, we define the set $O_{(c,s)}(\mathbf{d}) \triangleq \{i \in \mathcal{N} | d_i = (\cdot, c, s)\}$ of all WDs that use

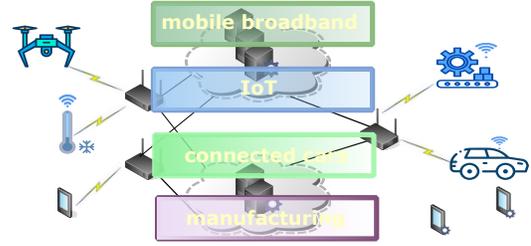


Fig. 1. An example of a slicing enabled MEC system that consists of $N = 7$ WDs, $C = 2$ ECs and $A = 3$ APs and $S = 4$ slices.

EC c in slice s and the set $O_c(\mathbf{d}) = \cup_{s \in \mathcal{S}} O_{(c,s)}(\mathbf{d})$ of all WDs that use EC c . Finally, we define the local computing singleton set $O_i(\mathbf{d}) \subset \{i, \emptyset\}$ for WD i (i.e., $O_i(\mathbf{d}) = \{i\}$ when WD i performs the computation locally and $O_i(\mathbf{d}) = \emptyset$ otherwise) and the set $O_l(\mathbf{d}) = \cup_{i \in \mathcal{N}} O_i(\mathbf{d})$ of all WDs that perform the computation locally.

Figure 1 shows an example of a slicing enabled MEC system that consists of $N = 7$ WDs, $C = 2$ ECs and $A = 3$ APs and $S = 4$ slices. In this example we have that 2 out of 7 WDs perform the computation locally and 5 out of 7 WDs offload their tasks. In what follows we discuss our models of communication and computing resources.

A. Communication Resources

Communication resources in the system are managed at two levels: at the network level and at the slice level.

At the network level, the radio resources of each AP $a \in \mathcal{A}$ are shared across the slices according to the *inter-slice radio resource allocation policy* $\mathcal{P}_b : \mathfrak{D} \rightarrow \mathbb{R}_{[0,1]}^{|\mathcal{A}| \times |\mathcal{S}|}$, which determines the inter-slice radio resource provisioning coefficients $b_a^s \in [0, 1]$, $\forall (a, s) \in \mathcal{A} \times \mathcal{S}$ such that $\sum_{s \in \mathcal{S}} b_a^s \leq 1$, $\forall a \in \mathcal{A}$.

At the slice level, the radio resources assigned to each slice $s \in \mathcal{S}$ are shared among the WDs according to an *intra-slice radio resource allocation policy* $\mathcal{P}_{w_a}^s : \mathfrak{D} \rightarrow \mathbb{R}_{[0,1]}^{|\mathcal{A}| \times |\mathcal{N}|}$, which determines the intra-slice radio resource provisioning coefficients $w_{i,a}^s \in [0, 1]$, $\forall a \in \mathcal{A}$ and $\forall i \in O_{(a,s)}(\mathbf{d})$ such that $\sum_{i \in O_{(a,s)}(\mathbf{d})} w_{i,a}^s \leq 1$, $\forall (a, s) \in \mathcal{A} \times \mathcal{S}$.

We denote by $R_{i,a}$ the achievable PHY rate of WD i at AP a . $R_{i,a}$ depends on physical signal characteristics, such as path loss and fading, and on the modulation-coding scheme. Given $R_{i,a}$ we can express the actual uplink rate of WD i at AP a in slice s as

$$W_{i,a}^s(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}^s) = b_a^s w_{i,a}^s R_{i,a}. \quad (1)$$

The uplink rate (1) together with the input data size D_i determines the transmission time of WD $i \in O_{(a,s)}(\mathbf{d})$,

$$T_{i,a}^{tx,s}(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}^s) = \frac{D_i}{W_{i,a}^s(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}^s)}. \quad (2)$$

Similar to previous works [10], [15], [16], [17] we make the assumption that the time needed to transmit the

results of the computation from the EC to the WD can be neglected because for many applications (e.g., face recognition and tracking) the size of the output data is significantly smaller than the size D_i of the input data.

B. Computing Resources

Our system model distinguishes between edge cloud resources and local computing resources.

1) *Edge Cloud Resources*: We consider that each slice $s \in \mathcal{S}$ is equipped with a certain combination of computing resources optimized for executing specific types of tasks (e.g, CPUs, GPUs, NPUs, FPGAs), and we denote by F_c^s the computing capability of EC c in slice s . The computing resources within a slice are shared among the WDs according to the *intra-slice computing power allocation policy* $\mathcal{P}_{w_c}^s : \mathcal{D} \rightarrow \mathbb{R}_{[0,1]}^{|\mathcal{C}| \times |\mathcal{N}|}$, which determines the intra-slice computing power provisioning coefficients $w_{i,c}^s \in [0, 1]$, $\forall c \in \mathcal{C}$ and $\forall i \in O_{(c,s)}(\mathbf{d})$ such that $\sum_{i \in O_{(c,s)}(\mathbf{d})} w_{i,c}^s = 1$, $\forall (c, s) \in \mathcal{C} \times \mathcal{S}$.

Given the computing capability F_c^s we can express the computing capability allocated to WD i in EC c in slice s as

$$F_{i,c}^s(\mathbf{d}, \mathcal{P}_{w_c}^s) = w_{i,c}^s F_c^s. \quad (3)$$

In order to account for the diversity of computing resources provided by different slices we use the coefficient $h_{i,s} \in \mathbb{R}_{\geq 0}$ to capture how well a slice s is tailored for executing a task generated by WD i and we express the expected number of instructions $L_{i,s}$ required to execute a task generated by WD i in slice s as $L_{i,s} = L_i/h_{i,s}$ (i.e., a high $h_{i,s}$ indicates that a task generated by WD i is a good match for the computing resources in slice s). Thus, in our model the computing capability (3) together with the expected task complexity $L_{i,s}$ determines the task execution time of WD $i \in O_{(c,s)}(\mathbf{d})$ as

$$T_{i,c}^{ex,s}(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}^s) = \frac{L_{i,s}}{F_{i,c}^s(\mathbf{d}, \mathcal{P}_{w_c}^s)}. \quad (4)$$

2) *Local Computing Resources*: We denote by F_i^l the computing capability of WD i and we express the local execution time T_i^{ex} of WD i as

$$T_i^{ex} = \frac{L_i}{F_i^l}. \quad (5)$$

C. Cost Model

We define the system cost as the aggregate completion time of all WDs. Before providing a formal definition, we introduce the shorthand notation

$$E_{i,e}^s = \begin{cases} \frac{D_i}{R_{i,e}^s} & \text{if } i \in \mathcal{N}, e \in \mathcal{E} \cap \mathcal{A}, s \in \mathcal{S} \\ \frac{L_{i,s}}{F_e^s} & \text{if } i \in \mathcal{N}, e \in \mathcal{E} \cap \mathcal{C}, s \in \mathcal{S}, \end{cases} \quad (6)$$

$$b_e^s = \begin{cases} b_e^s & \text{if } e \in \mathcal{E} \cap \mathcal{A}, s \in \mathcal{S} \\ 1 & \text{if } e \in \mathcal{E} \cap \mathcal{C}, s \in \mathcal{S}. \end{cases} \quad (7)$$

Cost of WD i : When offloading, the task completion time consists of two parts: the time needed to transmit

the data pertaining to a task through an AP and the time needed to execute a task in an EC. In the case of local computing, the task completion time depends only on the local execution time. Therefore, the cost of WD i can be expressed as

$$C_i(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}^s, \mathcal{P}_{w_c}^s) = \begin{cases} \frac{E_{i,a}^s}{b_a^s w_{i,a}^s} + \frac{E_{i,c}^s}{w_{i,c}^s}, & I_{\{d \neq (a,c,s)\}} = 1, \\ T_i^{ex}, & I_{\{d_i=i\}} = 1. \end{cases} \quad (8)$$

where $I_{\{d_i=d\}} = 1$ if $d_i = d$ and $I_{\{d_i=d\}} = 0$ otherwise.

Cost per slice: We express the cost in slice s as

$$C^s(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}^s, \mathcal{P}_{w_c}^s) = \sum_{e \in \mathcal{E}} \sum_{i \in O_{(e,s)}(\mathbf{d})} \frac{E_{i,e}^s}{b_e^s w_{i,e}^s}. \quad (9)$$

System cost: Finally, we express the system cost as

$$C(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}) = \sum_{s \in \mathcal{S}} C^s(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}^s, \mathcal{P}_{w_c}^s) + \sum_{i \in O_l(\mathbf{d})} C_i^l, \quad (10)$$

where $(\mathcal{P}_{w_a}, \mathcal{P}_{w_c}) = ((\mathcal{P}_{w_a}^s, \mathcal{P}_{w_c}^s))_{s \in \mathcal{S}}$ denotes the collection of slices' policies.

III. PROBLEM FORMULATION

We consider that the network operator aims at minimizing the system cost $C(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c})$ by finding an optimal vector $\hat{\mathbf{d}}$ of offloading decisions, and an optimal collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of policies for sharing the edge resources across slices and within slices. We refer to the problem as the *Joint Slice Selection and Edge Resource Management (JSS-ERM)* problem. Since the WDs generate atomic tasks that cannot be further split, the JSS-ERM is a mixed-integer optimization problem, and can be formulated as

$$\min_{\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}} C(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}) \quad (11)$$

$$\text{s.t. } \sum_{d \in \mathcal{D}_i} I_{\{d_i=d\}} = 1, \forall i \in \mathcal{N}, \quad (12)$$

$$C_i(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}) \leq T_i^{ex}, \forall i \in \mathcal{N}, \quad (13)$$

$$\sum_{s \in \mathcal{S}} b_a^s \leq 1, \forall a \in \mathcal{A}, \quad (14)$$

$$\sum_{j \in O_{(e,s)}(\mathbf{d})} w_{j,e}^s \leq 1, \forall e \in \mathcal{E}, \forall s \in \mathcal{S}, \quad (15)$$

$$b_a^s \geq 0, \forall a \in \mathcal{A}, \forall s \in \mathcal{S}, \quad (16)$$

$$w_{i,e}^s \geq 0, \forall i \in \mathcal{N}, \forall e \in \mathcal{E}, \forall s \in \mathcal{S}. \quad (17)$$

Constraint (12) enforces that each WD either performs the computation locally or offloads its task to exactly one logical resource $(a, c, s) \in \mathcal{A} \times \mathcal{C} \times \mathcal{S}$; constraint (13) ensures that the task completion time in the case of offloading is not greater than the task completion time in the case of local computing; constraint (14) enforces a limitation on the amount of communication resources of an AP that can be provided to each slice; constraint (15) enforces a limitation on the amount of communication resources of an AP and the amount of computing resources of an EC that can be provided to each WD in each slice.

Theorem 1. *The JSS-ERM defined by (11)-(17) is NP-hard.*

Proof. We provide the proof in Section IV-B. \square

In what follows we develop an approximation scheme for the JSS-ERM problem based on decomposition of the problem, and by adopting a game theoretic interpretation of one of the subproblems.

IV. NETWORK SLICE ORCHESTRATION AND EDGE RESOURCE ALLOCATION

In what follows we show that the JSS-ERM problem can be solved through solving a series of smaller optimization problems. To do so, we start with considering the problem of finding the collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal resource allocation policies for a given vector \mathbf{d} of offloading decisions.

Lemma 1. *Consider an offloading decision vector \mathbf{d} for which the constraint (13) can be satisfied. Furthermore, define the problem of finding a collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal resource allocation policies as*

$$\min_{\mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}} C(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}) \quad (18)$$

$$s.t. (13) - (17). \quad (19)$$

Then, the collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal resource allocation policies sets the provisioning coefficients according to

$$\hat{w}_{i,e}^s = \frac{\sqrt{E_{i,e}^s}}{\sum_{j \in O_{(e,s)}(\mathbf{d})} \sqrt{E_{j,e}^s}}, \forall e \in \mathcal{E}, \forall s \in \mathcal{S}, \forall i \in O_{(e,s)}(\mathbf{d}), \quad (20)$$

$$\hat{b}_a^s = \frac{\sum_{j \in O_{(a,s)}(\mathbf{d})} \sqrt{E_{j,a}^s}}{\sum_{s' \in \mathcal{S}} \sum_{j \in O_{(a,s')}(\mathbf{d})} \sqrt{E_{j,a}^{s'}}}, \forall a \in \mathcal{A}, \forall s \in \mathcal{S}. \quad (21)$$

Proof. First, observe that constraint (13) can be omitted since we assumed that the decision vector \mathbf{d} is such that constraint (13) can be satisfied. Furthermore, by inspecting the leading minors of the Hessian matrix of the objective function (18) it is easy to show that the matrix is positive semidefinite on the domain defined by (19), and thus problem (18)-(19) is convex. Therefore, the optimal solution of the problem must satisfy the Karush–Kuhn–Tucker (KKT) conditions and thus we can formulate the corresponding Lagrangian dual problem. To do so, let us define $\mathbf{b} \triangleq (b_a^s)_{a \in \mathcal{A}, s \in \mathcal{S}}$ and $\mathbf{w}^s \triangleq (w_{i,e}^s)_{i \in \mathcal{N}, e \in \mathcal{E}}$, and let us introduce non-negative Lagrange multiplier vectors $\boldsymbol{\alpha} = (\alpha_a)_{a \in \mathcal{A}}$, $\boldsymbol{\beta} = (\beta_e^s)_{e \in \mathcal{E}, s \in \mathcal{S}}$, $\boldsymbol{\gamma} = (\gamma_a^s)_{a \in \mathcal{A}, s \in \mathcal{S}}$ and $\boldsymbol{\delta} = (\delta_{i,e}^s)_{i \in O_{(e,s)}(\mathbf{d}), e \in \mathcal{E}, s \in \mathcal{S}}$ for constraints in (19), respectively. Next, let us define the Lagrangian dual problem corresponding to problem (18)-(19) as $\max_{\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\delta} \geq 0} \min_{\mathbf{b}, \mathbf{w} \geq 0} \mathcal{L}(\mathbf{b}, \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\delta})$, where the Lagrangian is given by

$$\begin{aligned} \mathcal{L}(\mathbf{b}, \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\delta}) = & \sum_{s' \in \mathcal{S}} \sum_{e' \in \mathcal{E}} \frac{1}{b_{e'}^{s'}} \left(\sum_{j \in O_{(e',s')}(\mathbf{d})} \frac{E_{j,e'}^{s'}}{w_{j,e'}^{s'}} \right) + \\ & \sum_{a' \in \mathcal{A}} \alpha_{a'} \left(\sum_{s' \in \mathcal{S}} b_{a'}^{s'} - 1 \right) + \sum_{e' \in \mathcal{E}} \sum_{s' \in \mathcal{S}} \beta_{e'}^{s'} \left(\sum_{j \in O_{(e',s')}(\mathbf{d})} w_{j,e'}^{s'} - 1 \right) \\ & - \sum_{a' \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \gamma_{a'}^{s'} b_{a'}^{s'} - \sum_{e' \in \mathcal{E}} \sum_{s' \in \mathcal{S}} \sum_{j \in O_{(e',s')}(\mathbf{d})} \delta_{j,e'}^{s'} w_{j,e'}^{s'} + \sum_{j \in O_t(\mathbf{d})} C_j^l. \end{aligned}$$

Now, we can express the KKT conditions as follows

$$\text{stationarity: } \sum_{j \in O_{(a,s)}(\mathbf{d})} \frac{E_{j,a}^s}{w_{j,a}^s (b_a^s)^2} = \alpha_a - \gamma_a^s, a \in \mathcal{A}, s \in \mathcal{S}, \quad (22)$$

$$\frac{E_{i,e}^s}{b_e^s (w_{i,e}^s)^2} = \beta_e^s - \delta_{i,e}^s, e \in \mathcal{E}, s \in \mathcal{S}, i \in O_{(e,s)}(\mathbf{d}), \quad (23)$$

$$\text{pr. feasibility: } (19), \quad (24)$$

$$\text{du. feasibility: } \boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\gamma}, \boldsymbol{\delta} \geq 0, \quad (25)$$

$$\text{co. slackness: } \alpha_a \left(\sum_{s' \in \mathcal{S}} b_{a'}^{s'} - 1 \right), a \in \mathcal{A}, \quad (26)$$

$$\beta_e^s \left(\sum_{j \in O_{(e,s)}(\mathbf{d})} w_{j,e}^s - 1 \right) = 0, e \in \mathcal{E}, s \in \mathcal{S}, \quad (27)$$

$$-\gamma_a^s b_a^s = 0, a \in \mathcal{A}, s \in \mathcal{S} \quad (28)$$

$$-\delta_{i,e}^s w_{i,e}^s = 0, e \in \mathcal{E}, s \in \mathcal{S}, i \in O_{(e,s)}(\mathbf{d}). \quad (29)$$

We proceed with finding $\hat{w}_{i,e}^s$. First, from the KKT dual feasibility condition $\boldsymbol{\delta} \geq 0$ and complementary slackness condition (29) we obtain that $\delta_{i,e}^s = 0$ must hold for every $e \in \mathcal{E}$, $s \in \mathcal{S}$ and $i \in O_{(e,s)}(\mathbf{d})$ as otherwise $w_{i,e}^s = 0$ would lead to infinite value of the objective function. Then, from the KKT stationarity condition (23) and complementary slackness condition (27) we obtain the expression (20) for coefficients $\hat{w}_{i,e}^s$. Finally, by substituting expression (20) into the KKT stationarity condition (22) and by following the same approach as for finding $\hat{w}_{i,e}^s$ we obtain the expression (21) for coefficients \hat{b}_a^s , which proves the result. \square

As a first step in the decomposition, let us consider the problem of finding the optimal collection $(\mathcal{P}_{w_a}^*, \mathcal{P}_{w_c}^*) = ((\mathcal{P}_{w_a}^{s,*}, \mathcal{P}_{w_c}^{s,*}))_{s \in \mathcal{S}}$ of resource allocation policies of slices for a given vector \mathbf{d} of offloading decisions and a given policy \mathcal{P}_b .

Proposition 1. *Consider an offloading decision vector \mathbf{d} for which constraint (13) can be satisfied and a policy \mathcal{P}_b for setting the inter-slice radio resource provisioning coefficients $b_a^s, \forall a \in \mathcal{A}, \forall s \in \mathcal{S}$. Then the solution to the problem*

$$\min_{\mathcal{P}_{w_a}^s, \mathcal{P}_{w_c}^s} C^s(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}^s, \mathcal{P}_{w_c}^s) \quad (30)$$

$$s.t. (13), (15), (17). \quad (31)$$

is given by (20), i.e., $(\mathcal{P}_{w_a}^{s,*}, \mathcal{P}_{w_c}^{s,*}) = (\hat{\mathcal{P}}_{w_a}^s, \hat{\mathcal{P}}_{w_c}^s), \forall s \in \mathcal{S}$.

Proof. The result can be proved by following the approach presented in the proof of Lemma 1. \square

As a second step, let us consider the problem of finding an optimal policy \mathcal{P}_b^* for a given vector \mathbf{d} of offloading decisions and the optimal collection $(\mathcal{P}_{w_a}^*, \mathcal{P}_{w_c}^*) = (\hat{\mathcal{P}}_{w_a}^*, \hat{\mathcal{P}}_{w_c}^*)$ of the slices' policies.

Proposition 2. Consider an offloading decision vector \mathbf{d} for which the constraint (13) can be satisfied. Furthermore, let us substitute (20) into (11)-(17) and define the problem of finding an optimal inter-slice radio resource allocation policy \mathcal{P}_b^* , i.e., a solution to

$$\min_{\mathcal{P}_b} \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} \frac{1}{b_{a'}^{s'}} \left(\sum_{j \in O_{(a', s')}}(\mathbf{d}) \sqrt{E_{j, a'}^{s'}} \right)^2 \quad (32)$$

$$\text{s.t. (13), (14) and (16)}. \quad (33)$$

Then, the optimal inter-slice radio resource allocation policy \mathcal{P}_b^* sets the inter-slice provisioning coefficients according to (21), i.e., $\mathcal{P}_b^* = \hat{\mathcal{P}}_b$.

Proof. The result can be proved by following the approach presented in the proof of Lemma 1. \square

By combining the above two results, we are now ready to show that the JSS-ERM problem can be decomposed into a sequence of optimization problems.

Theorem 2. The problem (18)-(19) can be solved optimally by finding the optimal policies $(\hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ first, and finding the optimal policy $\hat{\mathcal{P}}_b$ second, i.e.,

$$\begin{aligned} & \min_{\mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}} C(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}) = \\ & \min_{\mathcal{P}_b} \min_{\mathcal{P}_{w_a}, \mathcal{P}_{w_c}} C(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}) \end{aligned} \quad (34)$$

Proof. The result follows from the proofs of Lemma 1, Proposition 1 and Proposition 2. \square

Furthermore, as the next theorem shows, we can use this decomposition structure also for computing the optimal offloading decision vector.

Theorem 3. The problem (11)-(17) can be solved optimally by finding the optimal collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of resource allocation policies first, and finding an optimal offloading decision vector $\hat{\mathbf{d}}$ second, i.e.,

$$\begin{aligned} & \min_{\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}} C(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}) = \\ & \min_{\mathbf{d}} \min_{\mathcal{P}_b} \min_{\mathcal{P}_{w_a}, \mathcal{P}_{w_c}} C(\mathbf{d}, \mathcal{P}_b, \mathcal{P}_{w_a}, \mathcal{P}_{w_c}) \end{aligned} \quad (35)$$

Proof. It is easy to see that the exact values of the provisioning coefficients are functions of $\hat{\mathbf{d}}$. However, the optimal policies according to which the resources are shared are the same for every offloading decision vector $\mathbf{d} \in \mathcal{D}$, as defined by (20) and (21). Therefore, one can solve the problem (18)-(19) first, assuming an arbitrary offloading decision vector \mathbf{d} , and then given the solution $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of (18)-(19) find the optimal offloading decision vector $\hat{\mathbf{d}}$ that will determine the exact values of the provisioning coefficients. This proves the result. \square

A. Discussion and Practical Implications

So far we have shown that the JSS-ERM problem can be decomposed into a $S + 2$ coupled resource allocation

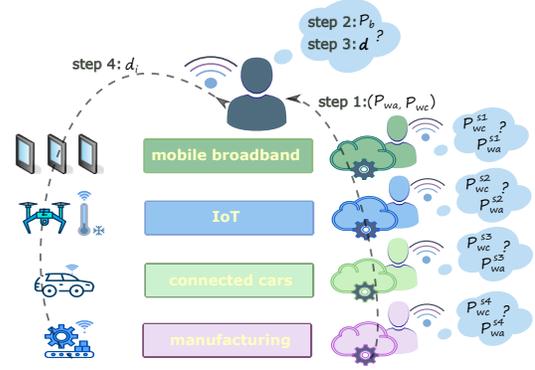


Fig. 2. An example of the potential implementation of a resource allocation and orchestration framework.

problems that can be solved sequentially. It is of interest to discuss the relationship between the decomposition and the potential implementation of a resource allocation and orchestration framework.

The proposed decomposition results in an optimization problem to be solved at the network level (eqns. (32)-(33)) and one in each slice ((eqns. (30)-(31))), followed by the problem of finding an optimal offloading decision vector. This structure is aligned with the slice-based network architecture proposed in [6], where inter-slice radio resource allocation and service orchestration are performed by a centralized entity, the *slice resource orchestrator* (SRO), while intra-slice radio and computing resource management is performed by the slices themselves, i.e., each slice manages its own radio and computing resources.

Figure 2 illustrates the interaction between the SRO and slices in the potential implementation of a resource allocation and orchestration framework.

B. Problem Complexity

In what follows we provide a result concerning the complexity of the JSS-ERM problem. For notational convenience let us first define the set of resources $\tilde{\mathcal{R}} \triangleq \{\{A \times \mathcal{S}\} \cup \{C \times \mathcal{S}\} \cup \mathcal{N}\}$ and let us introduce the following shorthand notation

$$\begin{aligned} q_{i,(a,s)} & \triangleq \sqrt{\frac{D_i}{R_{i,a}}}, \quad q_{i,(c,s)} \triangleq \sqrt{\frac{L_i}{h_{i,s}}}, \quad q_{i,i} \triangleq \sqrt{L_i}, \\ q_r(\mathbf{d}) & \triangleq \sum_{j \in O_r(\mathbf{d})} q_{j,r}, \quad \forall r \in \tilde{\mathcal{R}}. \end{aligned} \quad (36)$$

First, by substituting (20) into (8) and by using the notation introduced in (36), we can express the cost of WD i under a policy \mathcal{P}_b and the collection $(\hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal allocation policies of slices as

$$\tilde{C}_i(\mathbf{d}) = \sum_{r \in \tilde{\mathcal{R}}_{d_i}} m_r q_{i,r} q_r(\mathbf{d}), \quad (37)$$

where $\tilde{\mathcal{R}}_{d_i}$ is the set of resources that WD i uses for performing its task in \mathbf{d} (i.e., $\tilde{\mathcal{R}}_{d_i} \subset \tilde{\mathcal{R}}$) and $m_{(a,s)} = 1/b_a^s$, $m_{(c,s)} = 1/F_c^s$ and $m_i = 1/F_i^l$.

Second, by summing the expressions (37) over all WDs $i \in \mathcal{N}$ and by reordering the summations we can express

the system cost (10) under a policy \mathcal{P}_b and the collection $(\hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal allocation policies of slices as

$$\tilde{C}(\mathbf{d}) = \sum_{r \in \tilde{\mathcal{R}}} m_r q_r^2(\mathbf{d}). \quad (38)$$

Next, let us define the set of resources $\tilde{\mathcal{R}} \triangleq \{\mathcal{A} \cup \{\mathcal{C} \times \mathcal{S}\} \cup \mathcal{N}\}$ and a coefficient $q_{i,a} \triangleq q_{i,(a,s)} = \sqrt{D_i/R_{i,a}}$. By substituting (21) into (37), we can express the cost of WD i under the collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal allocation policies as

$$\tilde{C}_i(\mathbf{d}) = \sum_{r \in \tilde{\mathcal{R}}_{d_i}} m_r q_{i,r} q_r(\mathbf{d}), \quad (39)$$

where $\tilde{\mathcal{R}}_{d_i}$ is the set of resources that WD i uses for performing its task in \mathbf{d} (i.e., $\tilde{\mathcal{R}}_{d_i} \subset \tilde{\mathcal{R}}$) and $m_a = 1$.

Finally, by summing the expressions (39) over all WDs $i \in \mathcal{N}$ and by reordering the summations we can express the system cost (38) under the collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal allocation policies as

$$\bar{C}(\mathbf{d}) = \sum_{r \in \tilde{\mathcal{R}}} m_r q_r^2(\mathbf{d}). \quad (40)$$

Theorem 4. Consider the problem of finding the optimal vector $\hat{\mathbf{d}}$ of offloading decisions of WDs under the collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal allocation policies that set provisioning coefficients according to (20) and (21)

$$\min_{\mathbf{d}} \bar{C}(\mathbf{d}) \quad (41)$$

$$s.t. (12). \quad (42)$$

Problem (41)-(42) is NP-hard.

Proof. We prove the NP-hardness of the problem by reduction from the *Minimum Sum of Squares* problem (SP19 problem in [18]): given a finite set \mathcal{B} , a size $s(b) \in \mathbb{Z}^+$, $\forall b \in \mathcal{B}$ and positive integers $K \leq |\mathcal{B}|$ and J , the question is whether \mathcal{B} can be partitioned into K disjoint subsets $\mathcal{B}_1, \mathcal{B}_2, \dots, \mathcal{B}_K$ such that $\sum_{k=1}^K \left(\sum_{b \in \mathcal{B}_k} s(b) \right)^2 \leq J$.

For the reduction we set $S = 1$, $C = 0$ and $F_i^l = 0$, $\forall i \in \mathcal{N}$, i.e., in this simplified version of the problem $\tilde{\mathcal{R}} = \mathcal{A}$. Next, we let $\mathcal{N} = \mathcal{B}$, $|\mathcal{A}| = K$, $R_{i,a} = R_i$, $\forall i \in \mathcal{N}$, $\forall a \in \mathcal{A}$ and $\sqrt{D_i/R_i} = s(b)$. Then, it follows from (38) that the optimal solution of (41)-(42) provides the solution to the SP19 problem. As SP19 is NP-hard, problem (41)-(42) is also NP-hard, which proves the theorem. \square

Proof of Theorem 1. The result follows from Theorem 3 and Theorem 4. \square

V. APPROXIMATION SCHEME FOR THE JSS-ERM PROBLEM

In what follows we propose the *choose offloading slice* (COS) algorithm for computing an approximation to the optimal solution of the JSS-ERM problem. In particular, the algorithm serves as an approximation scheme to the problem of finding an optimal offloading decision vector.

$\mathbf{d}^* = \text{COS}(\mathbf{d}^0, \mathcal{P}_b, \mathcal{P}_{w_a}^*, \mathcal{P}_{w_c}^*)$

```

1  $\mathbf{d} \leftarrow \mathbf{d}^0$ 
2 while  $\exists \text{WD } j \in \mathcal{N}$  s.t.  $d_j \neq \arg \min_{d'_j \in \mathcal{D}_j} \tilde{C}_j(d'_j, d_{-j})$ 
3    $d_j^* = \arg \min_{d'_j \in \mathcal{D}_j} \tilde{C}_j(d'_j, d_{-j})$ ,  $\mathbf{d} = (d_j^*, d_{-j})$ 
4 end
5  $\mathbf{d}^* = \mathbf{d}$ 

```

Fig. 3. Pseudo code of the COS algorithm.

The algorithm starts from an offloading decision vector \mathbf{d}^0 in which all WDs perform computation locally and it lets WDs update their offloading decisions one at a time, based on their local cost function $\tilde{C}_i(\mathbf{d})$. We show the pseudo code of the algorithm in Figure 3.

Theorem 5. The COS algorithm terminates after a finite number of the iterations for any allocation policy \mathcal{P}_b and the collection $(\hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal allocation policies of slices.

Proof. The proof is based on a game theoretic treatment of the problem

$$\min_{\mathbf{d}} \bar{C}(\mathbf{d}) \quad (43)$$

$$s.t. (12), \quad (44)$$

in which the inter-slice radio resource provisioning coefficients are set according to an arbitrary policy \mathcal{P}_b and the intra-slice radio and computing power provisioning coefficients are set according to the optimal policies $\hat{\mathcal{P}}_{w_a}$ and $\hat{\mathcal{P}}_{w_c}$, respectively.

In what follows we show that the problem (43)-(44) can be interpreted as a congestion game $\Gamma(\mathcal{P}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c}) = \langle \mathcal{N}, (\mathcal{D}_i)_{i \in \mathcal{N}}, (\tilde{C}_i)_{i \in \mathcal{N}} \rangle$ with resource-dependent weights $q_{i,r}$, $i \in \mathcal{N}$, $r \in \tilde{\mathcal{R}}$, and the cost of WD i in the resulting game is given by (37). First, observe that $q_{i,r}$ can be interpreted as the weight that WD i contributes to the congestion when using resource $r \in \tilde{\mathcal{R}}$ and thus $q_r(\mathbf{d})$ can be interpreted as the total congestion on resource r in strategy profile \mathbf{d} . This in fact implies that the cost (37) of WD i in strategy profile \mathbf{d} depends on its own resource-dependent weights $q_{i,r}$ and on the total congestion $q_r(\mathbf{d})$ on the resources it uses. Therefore, it follows from [19] that the problem (43)-(44) can be interpreted as a congestion game $\Gamma(\mathcal{P}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ with resource dependent weights. Consequently, the COS algorithm terminates after a finite number of iterations iff the game $\Gamma(\mathcal{P}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ has a pure strategy Nash equilibrium. ¹

¹A pure strategy Nash equilibrium of a strategic game is a collection \mathbf{d}^* of decisions (called a strategy profile) for which $\tilde{C}_i(d_i^*, d_{-i}^*) \leq \tilde{C}_i(d_i, d_{-i}^*)$, $\forall d_i$, where d_i^* and d_{-i}^* are standard game theoretical notations for an improvement step of player i and for the collection of decisions (strategies) of all players other than i , respectively.

Since the cost $c_r(\mathbf{d}) \triangleq m_r q_r(\mathbf{d})$ of sharing every resource $r \in \mathcal{R}$ is an affine function of the congestion $q_r(\mathbf{d})$ on resource r , it follows from Theorem 4.2 in [19] that the game $\Gamma(\mathcal{P}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ has the exact potential function² given by

$$\Psi(\mathbf{d}) = \sum_{i \in \mathcal{N}} \sum_{r \in \mathcal{R}_{d_i}} q_{j,r} c_r^{\leq i}(\mathbf{d}), \quad (46)$$

where $c_r^{\leq i}(\mathbf{d}) = m_r q_r^{\leq i}(\mathbf{d})$ and $q_r^{\leq i}(\mathbf{d}) = \sum_{\{j \in \mathcal{O}_r(\mathbf{d}) | j \leq i\}} q_{i,r}$.

It is well known that in a finite strategic game that admits an exact potential all improvement paths³ are finite [20] and thus the existence of the exact potential function (46) allows us to use the COS algorithm for computing a pure strategy Nash equilibrium \mathbf{d}^* of the game $\Gamma(\mathcal{P}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$, which proves the result. \square

Theorem 6. *The COS algorithm terminates after a finite number of the iterations for the collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal allocation policies.*

Proof. By following the same approach as in the proof of Theorem 5, it is easy to show that given the collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal allocation policies, the problem (41)-(42) can be interpreted as a congestion game $\Gamma(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c}) = \langle \mathcal{N}, (\mathcal{D}_i)_{i \in \mathcal{N}}, (\tilde{C}_i)_{i \in \mathcal{N}} \rangle$ with resource-dependent weights $q_{i,r}$, $i \in \mathcal{N}$, $r \in \mathcal{R}$, and the cost of WD i in the resulting game is given by (39).

Since $m_a = 1, \forall a \in \mathcal{A}$, the cost $c_r(\mathbf{d}) \triangleq m_r q_r(\mathbf{d})$ of sharing every resource $r \in \mathcal{R}$ is an affine function of the congestion on resource r . Therefore, the game $\Gamma(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ is also an exact potential game, and thus the COS algorithm computes a pure strategy Nash equilibrium \mathbf{d}^* of the game $\Gamma(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$, which proves the result. \square

In general, the number of improvement steps can be exponential in a potential game, but as we show next the COS algorithm can compute an equilibrium \mathbf{d}^* of offloading decisions efficiently.

Theorem 7. *The COS algorithm terminates in $\mathcal{O}(n \frac{C^{min}}{C^{max}} \log \frac{\sum_{i \in \mathcal{N}} T_i^{ex}}{\Psi^{min}})$ iterations, where $n \geq 1$, C^{min} and C^{max} are system parameter dependent constants and Ψ^{min} is the minimum value of the potential function.*

Proof. First, let us define the minimum cost that WD i can achieve as $C_i^{min} \triangleq \min\{C_i^l, \min_{(a,c,s) \in \mathcal{A} \times \mathcal{C} \times \mathcal{S}} (D_i/R_{i,a} + L_{i,s}/F_c^s)\}$ and let $C^{min} \triangleq \min_{i \in \mathcal{N}} C_i^{min}$. Furthermore, let

²A function $\Psi : \times_i(\mathcal{D}_i) \rightarrow \mathbb{R}$ is an exact potential for a finite strategic game if for an arbitrary strategy profile (d_i, d_{-i}) and for any improvement step d_i^* the following holds:

$$\Psi(d_i, d_{-i}) - \Psi(d_i^*, d_{-i}) = \tilde{C}_i(d_i, d_{-i}) - \tilde{C}_i(d_i^*, d_{-i}). \quad (45)$$

³An improvement path is a sequence of strategy profiles in which one player at a time changes its strategy through performing an improvement step.

we define the maximum cost that WD i can achieve if it was the only WD in the system as $C_i^{max} \triangleq \max\{C_i^l, \min_{(a,c,s) \in \mathcal{A} \times \mathcal{C} \times \mathcal{S}} (D_i/R_{i,a} + L_{i,s}/F_c^s)\}$, and let $C^{max} = \max_{i \in \mathcal{N}} C_i^{max}$.

Consider now an iteration of the COS algorithm where the offloading decision of WD i is updated from d_i to d_i^* . We can then write

$$\begin{aligned} \Psi(d_i, d_{-i}) - \Psi(d_i^*, d_{-i}) &= \tilde{C}_i(d_i, d_{-i}) - \tilde{C}_i(d_i^*, d_{-i}) \\ &\geq -C^{max} \geq -\frac{C^{max}}{C^{min}} \Psi(d_i, d_{-i}), \end{aligned} \quad (47)$$

where the equality follows from the definition of the exact potential function (45), the first inequality follows from the fact that $\tilde{C}_i(d_i, d_{-i}) - \tilde{C}_i(d_i^*, d_{-i}) > 0$ since d_i^* is an improvement step of WD i and the last inequality follows from the fact that $\Psi(d_i, d_{-i}) \geq C^{min}$ for any vector \mathbf{d} of offloading decisions.

Therefore, from (47) we obtain $\Psi(d_i^*, d_{-i}) \leq (1 + \frac{C^{max}}{C^{min}})$, i.e., the COS algorithm decreases the potential function by at least a factor of $(1 + \frac{C^{max}}{C^{min}})$. Next, observe that from the definition of the constants C^{max} and C^{min} we have $\frac{C^{max}}{C^{min}} \geq 1$. Hence, since $(1+x)^{\frac{n}{x}} \leq e^n$ holds for $x, n \geq 1$, we obtain that after every $n \frac{C^{min}}{C^{max}}$ iterations of the COS algorithm $(1 + \frac{C^{max}}{C^{min}})^{n \frac{C^{min}}{C^{max}}} \leq e^n$, and thus every $n \frac{C^{min}}{C^{max}}$ iteration decreases the potential function by a constant factor (n can be chosen as a smallest positive constant for which $n \frac{C^{min}}{C^{max}} \geq 1$). Furthermore, since the COS algorithm starts from an offloading decision vector \mathbf{d}^0 in which all WDs perform computation locally, the potential function begins at the value $\Psi(\mathbf{d}^0) = \sum_{i \in \mathcal{N}} T_i^{ex}$ and cannot drop lower than Ψ^{min} . Therefore, the COS algorithm converges in $\mathcal{O}(n \frac{C^{min}}{C^{max}} \log \frac{\sum_{i \in \mathcal{N}} T_i^{ex}}{\Psi^{min}})$ iterations, which proves the result. \square

In what follows we address the efficiency of the COS algorithm in terms of the cost approximation ratio.

Theorem 8. *The COS algorithm is a 2.62-approximation algorithm for the optimization problem (43)-(44) in terms of the system cost, i.e., $\frac{\tilde{C}(\mathbf{d}^*)}{\tilde{C}(\hat{\mathbf{d}})} \leq 2.62$.*

Proof. Let us denote by $\mathcal{D}^* \subseteq \mathcal{D}$ the set of all vectors of offloading decisions that can be computed using the COS algorithm given any policy \mathcal{P}_b and the collection $(\hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of the optimal resource allocation policies of slices. Furthermore, let us consider a vector $\mathbf{d}^* \in \mathcal{D}^*$ and an arbitrary vector $\hat{\mathbf{d}} \in \mathcal{D}$ of offloading decisions. Since there is no WD i for which the cost $\tilde{C}_i(\mathbf{d}^*)$ can be decreased by unilaterally changing its offloading decision we have the following

$$\tilde{C}_i(\mathbf{d}^*) \leq \sum_{r \in \tilde{\mathcal{R}}_{d_i^*} \cap \tilde{\mathcal{R}}_{\hat{d}_i}} m_r q_{i,r} q_r(\mathbf{d}^*) + \quad (48)$$

$$\sum_{r \in \tilde{\mathcal{R}}_{d_i^*} \setminus \tilde{\mathcal{R}}_{\hat{d}_i}} m_r (q_r(\mathbf{d}^*) + q_{i,r}) q_{i,r} \leq \sum_{r \in \tilde{\mathcal{R}}_{\hat{d}_i}} m_r (q_r(\mathbf{d}^*) + q_{i,r}) q_{i,r},$$

where $\tilde{\mathcal{R}}_{d_i^*} \subset \tilde{\mathcal{R}}$ and $\tilde{\mathcal{R}}_{\hat{d}_i} \subset \tilde{\mathcal{R}}$ denote the the set of resources that WD i uses in \mathbf{d}^* and $\hat{\mathbf{d}}$, respectively. By

summing (48) over all WDs $i \in \mathcal{N}$ and by reordering the summations we obtain

$$\tilde{C}(\mathbf{d}^*) \leq \sum_{r \in \mathcal{R}} \sum_{i \in \mathcal{O}_r(\hat{\mathbf{d}})} m_r (q_r(\mathbf{d}^*) q_{i,r} + q_{i,r}^2). \quad (49)$$

From the definition (36) of the total weight $q_r(\mathbf{d})$ on resource $r \in \tilde{\mathcal{R}}$ and from $\sum_{i \in \mathcal{O}_r(\hat{\mathbf{d}})} q_{i,r}^2 \leq q_r^2(\hat{\mathbf{d}})$ we obtain

$$\tilde{C}(\mathbf{d}^*) \leq \sum_{r \in \mathcal{R}} m_r q_r(\mathbf{d}^*) q_r(\hat{\mathbf{d}}) + \sum_{r \in \mathcal{R}} m_r q_r^2(\hat{\mathbf{d}}).$$

Next, let us recall the Cauchy-Schwartz inequality $\sum_{r \in \mathcal{R}} a_r b_r \leq \sqrt{\sum_{r \in \mathcal{R}} a_r^2 \sum_{r \in \mathcal{R}} b_r^2}$. By defining $a_r \triangleq \sqrt{m_r} q_r(\mathbf{d}^*)$ and $b_r \triangleq \sqrt{m_r} q_r(\hat{\mathbf{d}})$ we obtain the following

$$\tilde{C}(\mathbf{d}^*) \leq \sqrt{\sum_{r \in \mathcal{R}} m_r q_r^2(\mathbf{d}^*) \sum_{r \in \mathcal{R}} m_r q_r^2(\hat{\mathbf{d}})} + \sum_{r \in \mathcal{R}} m_r q_r^2(\hat{\mathbf{d}}). \quad (50)$$

By dividing the right and the left side of (50) by $\sum_{r \in \mathcal{R}} q_r^2(\hat{\mathbf{d}}) > 0$ and by using (38) we obtain

$$\frac{\tilde{C}(\mathbf{d}^*)}{\tilde{C}(\hat{\mathbf{d}})} \leq \sqrt{\frac{\tilde{C}(\mathbf{d}^*)}{\tilde{C}(\hat{\mathbf{d}})}} + 1. \quad (51)$$

Since (51) holds for any vector $\mathbf{d}^* \in \mathcal{D}^*$ of offloading decisions computed by the COS algorithm and for any vector $\hat{\mathbf{d}} \in \mathcal{D}$ of offloading decisions of the WDs, it holds for the worst vector $\mathbf{d}^* = \arg \max_{\mathbf{d} \in \mathcal{D}^*} \tilde{C}(\mathbf{d})$ of offloading decisions that can be computed using the COS algorithm and for the optimal $\hat{\mathbf{d}} = \arg \min_{\mathbf{d} \in \mathcal{D}} \tilde{C}(\mathbf{d})$ solution too. Therefore, by solving (51) we obtain that the cost approximation ratio $\frac{\tilde{C}(\mathbf{d}^*)}{\tilde{C}(\hat{\mathbf{d}})}$ of the COS algorithm is upper bounded by $(3 + \sqrt{5})/2 \cong 2.62$, which proves the theorem. \square

Theorem 9. *The COS algorithm is a 2.62-approximation algorithm for the optimization problem (41)-(42) in terms of the system cost, i.e., $\frac{\tilde{C}(\mathbf{d}^*)}{\tilde{C}(\hat{\mathbf{d}})} \leq 2.62$.*

Proof. The result can be easily obtained by following the approach used to prove Theorem 8. \square

Finally, from Theorem 3 and Theorem 9 we obtain the approximation ratio bound for the proposed decomposition-based algorithm.

Theorem 10. *Given the collection $(\hat{\mathcal{P}}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})$ of optimal allocation policies, the proposed decomposition-based algorithm computes a 2.62-approximation solution to the JSS-ERM problem.*

VI. NUMERICAL RESULTS

We used extensive simulations to evaluate the performance of the proposed resource allocation algorithm. To

capture the potentially uneven spatial distribution of ECs, WDs and APs in a dense urban area, we consider a square area of $1km \times 1km$ in which WDs and 3 ECs are placed uniformly at random and 5 APs are placed at random on a *regular grid* with 25 points. The channel gain of WD i to AP a depends on their Euclidean distance $d_{i,a}$ and on the path loss exponent α , which we set to 4 according to the path loss model in urban and suburban areas [21]. We set the bandwidth B_a of 2 APs to 18MHz and the bandwidth of 3 APs to 27MHz, corresponding to 25 and 75 resource blocks that are $12 \times 60KHz$ and $12 \times 30KHz$ subcarriers wide [22], [23], respectively. We consider that the transmit power $P_{i,a}$ of every WD i is uniformly distributed on $[10^{-6}, 0.1]W$ according to [24]. We calculate the total thermal noise in a $B_a MHz$ channel as $N_0(\text{dBm}) = -174 + 10 \log(B_a)$ according to [25] and the transmission rate $R_{i,a}$ achievable to WD i at AP a as $R_{i,a} = B_a \log(1 + d_{i,a}^{-\alpha} \frac{P_{i,a}}{N_0})$.

To set the values for the computational capabilities of the WDs, we consider a line of Samsung Galaxy phones, from the oldest version with 1 core operating at 1GHz to the one of the newest versions with 8 cores operating at 2.84GHz. We consider that EC c_1 is equipped with 36 vCPUs operating at 2.3GHz and 96 vCPUs operating at 3.6GHz. We consider that EC c_2 and EC c_3 are equipped with 1 GPU each (with 2048 parallel processing cores operating at 557MHz and 2496 parallel processing cores operating at 560MHz, respectively). Given the measurements reported in [26], [27], [28] we assume that a WD, a CPU and a GPU can execute on average 2, 3 and 1 instructions per cycle (IPC), respectively. Based on this, we consider that the computational capability F_i^l of every WD i is uniformly distributed on $[2, 45.4]GIPS$, where the lower and the upper bound correspond to the oldest and the newest version of the phone, respectively. Similarly, we calculate the computational capabilities of ECs, and set them to $F_{c_1} = 1285.2GIPS$, $F_{c_2} = 1140.7GIPS$ and $F_{c_3} = 1397.8GIPS$.

The input data size D_i is drawn from a uniform distribution on $[1.7, 10]Mb$ according to measurements in [29]. The number X of instructions per data bit follows a Gamma distribution [30] with shape parameter $k = 75$ and scale $\theta = 50$. Given D_i and X , we calculate the complexity of a task as $L_i = D_i X$.

Motivated by Amazon EC2 instances [31] designed to support different kinds of applications (e.g., G3 and P2 instances for graphics-intensive and general-purpose GPU applications, and C5 and I3 instances for compute-intensive and non-virtualized workloads), we evaluate the system performance for the following four cases.

S= 1: The slice s_1 contains all ECs, and thus is able to support all of the above applications.

S= 2: The ECs are sliced such that slice s_1 supports the G3.4 instance and slice s_2 supports instances C5 and I3.

S= 3: The ECs are sliced such that slices s_1 and s_2 support P2 and G3s instances, respectively and slice s_3

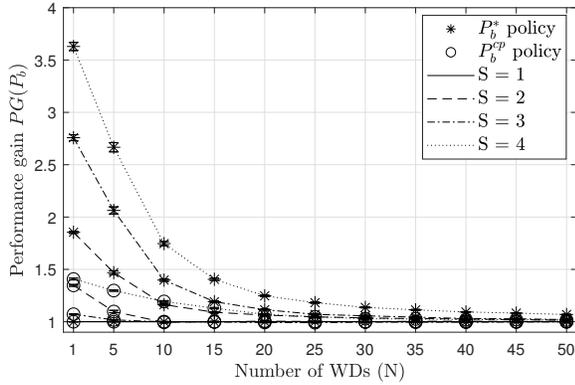


Fig. 4. Performance gain vs. number of WDs N .

supports instances C5 and I3.

S = 4: The ECs are sliced such that slices s_1, s_2, s_3 and s_4 support P2, G3s, C5 and I3 instances, respectively.

The coefficients $\frac{1}{h_{i,s}}$ were drawn from a continuous uniform distribution on $[0, 1]$ and unless otherwise noted, the results are shown for all of the above scenarios.

We use two bandwidth allocation policies \mathcal{P}_b of the slice orchestrator as a basis for comparison. The first policy \mathcal{P}_b^{cp} shares the bandwidth of each AP a among slices proportionally to the ECs' resources that slices have. The second policy \mathcal{P}_b^{eq} gives an equal share of the bandwidth of each AP a to each slice s . Observe that the COS algorithm computes an approximation vector \mathbf{d}^* of offloading decisions for both policies (c.f. Theorem 5 and Theorem 8). The results shown are the averages of 300 simulations, together with 95% confidence intervals.

A. System Performance

We start with considering the system performance from the point of view of the slice orchestrator. To do so we define the system performance gain $PG(\mathcal{P}_b)$ for an inter-slice radio allocation policy \mathcal{P}_b w.r.t. the policy \mathcal{P}_b^{eq} as

$$PG(\mathcal{P}_b) = \frac{C(\mathbf{d}^*, \mathcal{P}_b^{eq}, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})}{C(\mathbf{d}^*, \mathcal{P}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})}.$$

Figure 4 shows $PG(\mathcal{P}_b)$ as a function of the number N of WDs for the optimal \mathcal{P}_b^* and for the cloud proportional \mathcal{P}_b^{cp} allocation policy of the operator. We observe that $PG(\mathcal{P}_b^*) = PG(\mathcal{P}_b^{cp}) = 1$ when $S = 1$ because the three solutions are equivalent when there is no slicing. On the contrary, for $S > 1$ we observe that $PG(\mathcal{P}_b^*) > 1$ and $PG(\mathcal{P}_b^{cp}) > 1$, which is due to that the policy \mathcal{P}_b^{cp} does not take into account that the slices might have different amounts of ECs' resources. We also observe that the policy \mathcal{P}_b^* achieves better performance gain (up to 2.5 times greater) than the policy \mathcal{P}_b^{cp} because \mathcal{P}_b^* assigns the WDs to slices not only based on the amount of ECs' resources the slices have, but also based on how well the slices are tailored for executing their tasks. This effect is especially evident when there are few WDs, because in

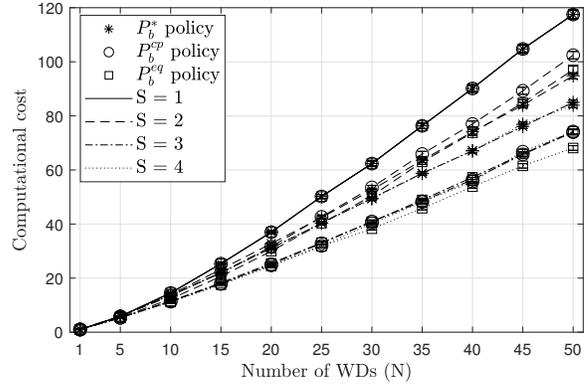


Fig. 5. Computational complexity vs. number of WDs N .

this case WDs tend to offload their tasks, and thus the system cost is mostly determined by the offloading cost.

B. Computational Cost

Figure 5 shows the number of iterations in which the COS algorithm computes a decision vector \mathbf{d}^* as a function of the number N of WDs under the optimal \mathcal{P}_b^* , the cloud proportional \mathcal{P}_b^{cp} and the equal \mathcal{P}_b^{eq} inter-slice radio allocation policy of the slice orchestrator.

Interestingly, the number of updates decreases with the number S of slices. This is due to that the congestion on the logical resources decreases as S increases, and thus the COS algorithm updates the offloading decisions less frequently. We also observe that the number of updates scales approximately linearly with N under all considered policies of the slice orchestrator, and thus we can conclude that the COS algorithm is computationally efficient, which makes it a good candidate for computing an approximation \mathbf{d}^* to the optimal vector $\hat{\mathbf{d}}$ of offloading decisions of WDs.

C. Performance Within the Slices

We continue with considering the performance from the point of view of the slices. For an inter-slice radio allocation policy \mathcal{P}_b , we denote by $n^s(\mathcal{P}_b)$ the number of offloaders per slice in the vector \mathbf{d}^* of offloading decisions computed by the COS algorithm and we define the cost ratio $CR^s(\mathcal{P}_b)$ per slice w.r.t. the system cost as

$$CR^s(\mathcal{P}_b) = \frac{C^s(\mathbf{d}^*, \mathcal{P}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})}{C(\mathbf{d}^*, \mathcal{P}_b, \hat{\mathcal{P}}_{w_a}, \hat{\mathcal{P}}_{w_c})}.$$

Figure 6 and Figure 7 show $n^s(\mathcal{P}_b)$ and $CR^s(\mathcal{P}_b)$, respectively for the optimal \mathcal{P}_b^* , the cloud proportional \mathcal{P}_b^{cp} and the equal \mathcal{P}_b^{eq} inter-slice radio allocation policy of the slice orchestrator. The results are shown for $S = 2$ and the red lines in Figure 7 show the share of the ECs' resources among the slices s_1 and s_2 (i.e. slices s_1 and s_2 have approximately 72% and 28% of the resources, respectively). We observe from Figure 6 and Figure 7, respectively that the gap between $n^{s_1}(\mathcal{P}_b)$ and $n^{s_2}(\mathcal{P}_b)$ and the gap between $CR^{s_1}(\mathcal{P}_b)$ and $CR^{s_2}(\mathcal{P}_b)$

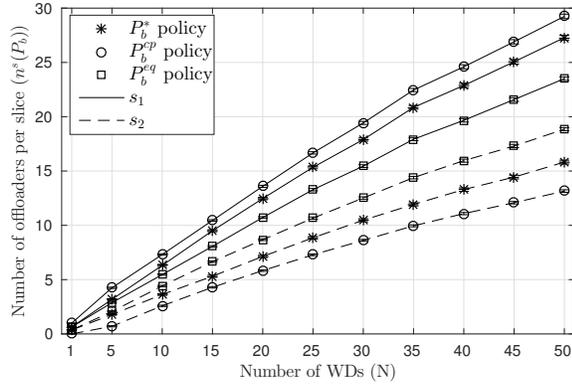


Fig. 6. Number of offloaders per slice vs. number of WDs N .

are highest in the case of the policy \mathcal{P}_b^{cp} and lowest in the case of the policy \mathcal{P}_b^{cq} . Therefore, WDs whose tasks are a better match with the EC resources in slice s_2 than those in slice s_1 cannot fully exploit the ECs' resources in slice s_2 under the policy \mathcal{P}_b^{cp} , which allocates bandwidth resources proportionally to the ECs' resources. Similarly, WDs whose tasks are a better match with the EC resources in slice s_1 than in slice s_2 cannot fully exploit the ECs' resources in slice s_1 under the policy \mathcal{P}_b^{cq} , which allocates bandwidth resources equally. On the contrary, the results show that the optimal policy \mathcal{P}_b^* finds a good match between the EC resources in the slices and the WDs' preferences for different types of computing resources, which makes it a good candidate for dynamic resource management for network slicing coupled with edge computing.

VII. RELATED WORK

Closest to our work a recent game theoretic treatments of the computation offloading problem [32], [33], [34], [35], [36]. In [32] the authors considered devices that compete for cloud resources so as to minimize their energy consumption, and proved that an equilibrium of offloading decision can be computed in polynomial time. In [33] the authors considered devices that maximize their performance and a profit maximizing service provider, and used backward induction for deriving near optimal strategies for the devices and the operator. In [34] the authors considered that devices can offload their tasks to a cloud through multiple identical wireless links, modeled the congestion on wireless links, and used a potential function argument for proposing a decentralized algorithm for computing an equilibrium. In [35] the authors considered that devices can offload their tasks to a cloud through multiple heterogeneous wireless links, modeled the congestion on wireless and cloud resources, showed that the game played by devices is not a potential game and proposed a decentralized algorithm for computing an equilibrium. In [36] the authors modeled the interaction between devices and a single network operator as a Stackelberg game, and provided an algorithm for computing

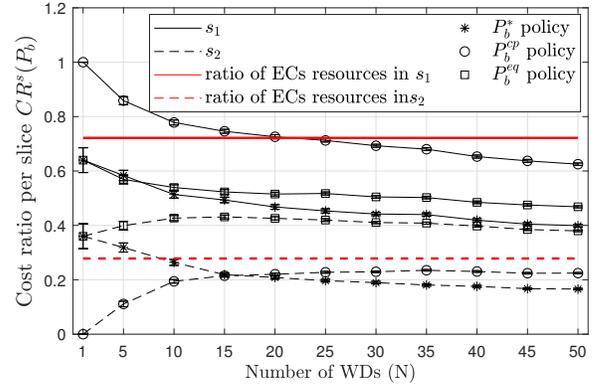


Fig. 7. Cost ratio per slice vs. number of WDs N .

a subgame perfect equilibrium. Unlike these works, we consider the computation offloading problem together with network slicing and we analyze the interaction between the network operator and the slices.

Another line of works considers the network slicing resource allocation problem [37], [38], [39], [40], [41]. In [37] the authors considered an auction-based model for allocating edge cloud resources to slices and proposed an algorithm for allocating resources to slices so as to maximize the total network revenue. In [38] the authors considered the radio resources slicing problem and proposed an approximation algorithm for maximizing the sum of the users' utilities. In [39] the authors modeled the interaction between slices that compete for bandwidth resources with the objective to maximize the sum of their users' utilities, and proposed an admission control algorithm under which the slices can reach an equilibrium. In [40] the authors proposed a deep learning architecture for sharing the resources among network slices in order to meet the users' demand within the slices. In [41] the authors considered a radio access network slicing problem and proposed two approximation algorithms for maximizing the total network throughput. Unlike these works, we consider a slicing enabled edge system in which the slice resource orchestrator assigns WDs to slices and shares radio resources across slices, while the slices manage their own radio and computing resources with the objective to maximize overall system performance.

To the best of our knowledge ours is the first work to consider slicing and computation offloading to edge clouds jointly, capturing the interaction between the slice resource orchestrator and the slices.

VIII. CONCLUSION

We have considered the computation offloading problem in an edge computing system under network slicing in which slices jointly manage their own communication and computing resources and the slice resource orchestrator manages communication resources among slices and assigns the WDs to slices. We formulated the problem of minimizing the sum over all WDs' task completion times

as a mixed-integer problem, proved that the problem is NP-hard and proposed a decomposition of the problem into a sequence of optimization problems. We proved that the proposed decomposition does not change the optimal solution of the original problem, proposed an efficient approximation algorithm for solving the decomposed problem and proved that the algorithm has bounded approximation ratio. Our numerical results show that the proposed algorithm is computationally efficient. They also show that dynamic allocation of slice resources is essential for maximizing the benefits of edge computing, and slicing could be beneficial for improving overall system performance.

REFERENCES

- [1] J. Ordonez-Lucena, P. Ameigeiras, D. Lopez, J. J. Ramos-Munoz, J. Lorca, and J. Folgueira, "Network slicing for 5g with sdn/nfv: Concepts, architectures, and challenges," *IEEE Communications Magazine*, vol. 55, no. 5, pp. 80–87, 2017.
- [2] S. Kekki, W. Featherstone, Y. Fang, P. Kuure, A. Li, A. Ranjan, D. Purkayastha, F. Jiangping, D. Frydman, G. Verin *et al.*, "Mec in 5g networks," *Sophia Antipolis, France, ETSI, White Paper*, 2018.
- [3] M. Rost and S. Schmid, "Virtual network embedding approximations: Leveraging randomized rounding," *IEEE/ACM Trans. Netw.*, vol. 27, no. 5, pp. 2071–2084, Oct. 2019.
- [4] B. Farkiani, B. Bakhshi, and S. A. MirHassani, "A fast near-optimal approach for energy-aware sfc deployment," *IEEE Transactions on Network and Service Management*, vol. 16, no. 4, pp. 1360–1373, Dec 2019.
- [5] I. Jang, D. Suh, S. Pack, and G. Dán, "Joint optimization of service function placement and flow distribution for service function chaining," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 11, pp. 2532–2541, Nov 2017.
- [6] A. Zafeiropoulos *et al.*, *5G PPP Architecture Working Group: View on 5G Architecture*, S. Redana and Ö. Bulakci, Eds. European Commission, Jun. 2019, vol. Version 3.0.
- [7] X. Foukas, M. K. Marina, and K. Kontovasilis, "Orion: Ran slicing for a flexible and cost-effective multi-service mobile network architecture," in *Proc. of ACM International Conference on Mobile Computing and Networking (MobiCom)*, 2017, pp. 127–140.
- [8] C. Chang, N. Nikaein, and T. Spyropoulos, "Radio access network resource slicing for flexible service execution," in *Proc. of IEEE INFOCOM 2018 Workshops*, April 2018, pp. 668–673.
- [9] J. Zheng, Y. Cai, Y. Wu, and X. Shen, "Dynamic computation offloading for mobile cloud computing: A stochastic game-theoretic approach," *IEEE Transactions on Mobile Computing*, vol. 18, no. 4, pp. 771–786, 2018.
- [10] X. Chen, "Decentralized computation offloading game for mobile cloud computing," *IEEE Transactions on Parallel and Distributed Systems*, vol. 26, no. 4, pp. 974–983, 2014.
- [11] S. Jošilo and G. Dán, "Decentralized algorithm for randomized task allocation in fog computing systems," *IEEE/ACM Transactions on Networking*, vol. 27, no. 1, pp. 85–97, 2018.
- [12] J. L. D. Neto, S.-Y. Yu, D. F. Macedo, J. M. S. Nogueira, R. Langar, and S. Secci, "Uloof: a user level online offloading framework for mobile edge computing," *IEEE Transactions on Mobile Computing*, vol. 17, no. 11, pp. 2660–2674, 2018.
- [13] B.-G. Chun, S. Ihm, P. Maniatis, M. Naik, and A. Patti, "Clonecloud: elastic execution between mobile device and cloud," in *Proceedings of the sixth conference on Computer systems*. ACM, 2011, pp. 301–314.
- [14] E. Cuervo, A. Balasubramanian, D.-k. Cho, A. Wolman, S. Saroiu, R. Chandra, and P. Bahl, "Maui: making smartphones last longer with code offload," in *Proceedings of the 8th international conference on Mobile systems, applications, and services*. ACM, 2010, pp. 49–62.
- [15] S. Jošilo and G. Dán, "Selfish decentralized computation offloading for mobile cloud computing in dense wireless networks," *IEEE TMC*, vol. 18, no. 1, pp. 207–220, 2018.
- [16] D. Huang, P. Wang, and D. Niyato, "A dynamic offloading algorithm for mobile computing," *IEEE Transactions on Wireless Communications*, vol. 11, no. 6, pp. 1991–1995, 2012.
- [17] S. Jošilo and G. Dán, "Joint management of wireless and computing resources for computation offloading in mobile edge clouds," *IEEE Transactions on Cloud Computing*, pp. 1–1, 2019.
- [18] M. R. Garey and D. S. Johnson, *Computers and intractability*. wh freeman New York, 2002, vol. 29.
- [19] T. Harks, M. Klimm, and R. H. Möhring, "Characterizing the existence of potential functions in weighted congestion games," *Theory of Computing Systems*, pp. 46–70, 2011.
- [20] D. Monderer and L. S. Shapley, "Potential games," *Games and economic behavior*, vol. 14, no. 1, pp. 124–143, 1996.
- [21] S. R. Saunders and A. Aragón-Zavala, *Antennas and propagation for wireless communication systems*. John Wiley & Sons, 2007.
- [22] E. TSGR, "Lte: Evolved universal terrestrial radio access (e-utra)," *Physical channels and modulation (3GPP TS 36.211 version 10.0.0. 0 Release 10) ETSI TS*, 2011.
- [23] A. Zaidi, F. Athley, J. Medbo, U. Gustavsson, G. Durisi, and X. Chen, *5G Physical Layer: Principles, Models and Technology Components*. Academic Press, 2018.
- [24] M. Lauridsen, L. Noël, T. B. Sørensen, and P. Mogensen, "An empirical lte smartphone power model with a view to energy efficiency evolution." *Intel Technology Journal*, vol. 18, no. 1, 2014.
- [25] N. Da Dalt and A. Sheikholeslami, *Understanding Jitter and Phase Noise: A Circuits and Systems Perspective*. Cambridge University Press, 2018.
- [26] L. Codrescu, W. Anderson, S. Venkumanhanti, M. Zeng, E. Plondke, C. Koob, A. Ingle, C. Tabony, and R. Maule, "Hexagon dsp: An architecture optimized for mobile multimedia and communications," *IEEE Micro*, vol. 34, no. 2, pp. 34–43, 2014.
- [27] D. Hackenberg, R. Schöne, T. Ilsche, D. Molka, J. Schuchart, and R. Geyer, "An energy efficiency feature survey of the intel haswell processor," in *2015 IEEE international parallel and distributed processing symposium workshop*, 2015, pp. 896–904.
- [28] Y. Takefuji, *GPU Parallel Computing for Machine Learning in Python: How to Build a Parallel Computer*. Independently published, 2017.
- [29] L. Fletcher, L. Petersson, and A. Zelinsky, "Road scene monotony detection in a fatigue management driver assistance system," in *IEEE Proceedings. Intelligent Vehicles Symposium, 2005.*, 2005, pp. 484–489.
- [30] J. R. Lorch and A. J. Smith, "Pace: A new approach to dynamic voltage scaling," *IEEE Transactions on Computers*, vol. 53, no. 7, pp. 856–869, 2004.
- [31] <https://aws.amazon.com/ec2/instance-types/>, "Amazon ec2 instance types."
- [32] Y. Ge, Y. Zhang, Q. Qiu, and Y.-H. Lu, "A game theoretic resource allocation for overall energy minimization in mobile cloud computing system," in *ACM/IEEE Symposium on low power electronics and design*, 2012, pp. 279–284.
- [33] Y. Wang, X. Lin, and M. Pedram, "A nested two stage game-based optimization framework in mobile cloud computing system," in *IEEE Service Oriented System Engineering Symposium*, 2013, pp. 494–502.
- [34] X. Chen, L. Jiao, W. Li, and X. Fu, "Efficient multi-user computation offloading for mobile-edge cloud computing," *IEEE/ACM Transactions on Networking*, vol. 24, no. 5, pp. 2795–2808, 2015.
- [35] S. Jošilo and G. Dán, "A game theoretic analysis of selfish mobile computation offloading," in *IEEE INFOCOM*, 2017, pp. 1–9.
- [36] —, "Wireless and computing resource allocation for selfish computation offloading in edge computing," in *IEEE INFOCOM*, 2019, pp. 2467–2475.
- [37] M. Jiang, M. Condoluci, and T. Mahmoodi, "Network slicing in 5g: An auction-based model," in *IEEE International Conference on Communications (ICC)*, 2017, pp. 1–6.
- [38] P. Caballero, A. Banchs, G. De Veciana, and X. Costa-Pérez, "Multi-tenant radio access network slicing: Statistical multiplex-

- ing of spatial loads,” *IEEE/ACM Transactions on Networking*, vol. 25, no. 5, pp. 3044–3058, 2017.
- [39] P. Caballero, A. Banchs, G. De Veciana, X. Costa-Pérez, and A. Azcorra, “Network slicing for guaranteed rate services: Admission control and resource allocation games,” *IEEE Transactions on Wireless Communications*, vol. 17, no. 10, pp. 6419–6432, 2018.
- [40] D. Bega, M. Gramaglia, M. Fiore, A. Banchs, and X. Costa-Pérez, “Deepcog: Cognitive network management in sliced 5g networks with deep learning,” in *IEEE INFOCOM*, 2019, pp. 280–288.
- [41] S. D’Oro, F. Restuccia, A. Talamonti, and T. Melodia, “The slice is served: Enforcing radio access network slicing in virtualized 5g systems,” in *IEEE INFOCOM*, 2019, pp. 442–450.