

Joint VNF Placement and CPU Allocation in 5G

Original

Joint VNF Placement and CPU Allocation in 5G / Agarwal, Satyam; Malandrino, Francesco; Chiasserini, Carla Fabiana; De, Swades. - STAMPA. - (2018). ((Intervento presentato al convegno IEEE INFOCOM 2018 tenutosi a Hawaii (USA) nel April 2018.

Availability:

This version is available at: 11583/2694650 since: 2018-03-07T09:02:58Z

Publisher:

IEEE

Published

DOI:

Terms of use:

openAccess

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)

Joint VNF Placement and CPU Allocation in 5G

S. Agarwal, F. Malandrino, C. F. Chiasserini, S. De

Abstract—Thanks to *network slicing*, 5G networks will support a variety of services in a flexible and swift manner. In this context, we seek to make high-quality, *joint* optimal decisions concerning the placement of VNFs across the physical hosts for realizing the services, and the allocation of CPU resources in VNFs sharing a host. To this end, we present a queuing-based system model, accounting for all the entities involved in 5G networks. Then, we propose a fast and efficient solution strategy yielding near-optimal decisions. We evaluate our approach in multiple scenarios that well represent real-world services, and find it to consistently outperform state-of-the-art alternatives and closely match the optimum.

I. INTRODUCTION

Enabled by network function virtualization (NFV) and software defined networking (SDN), *network slicing* [1] is one of the most exciting features of 5G networks. Third parties (“verticals”) will specify the services they want to provide and the associated requirements, e.g., maximum latency or minimum throughput, to the network operator. Thanks to NFV, such services will be expressed as graphs of virtual network functions (VNFs), running on virtual machines or containers. Through SDN, the VNF graph will then be mapped onto the physical and virtual resources available in the network, which can be seen as a pool of resources the operator can draw from.

In this scenario, there are three main entities to account for. First, *VNFs* (e.g., firewall or transcoding) performing the processing required by different types of services and running into virtual machines or containers; second, physical *hosts*, capable of running VNFs; third, the *links* connecting the physical hosts together. Additionally, there are two main decisions we need to make in order to effectively manage the network: (i) *VNF placement*, i.e., which VNFs run at each physical host, and (ii) *allocation*, i.e., how the computational capabilities available at physical hosts are allocated to the VNFs they run. VNF placement and CPU allocation¹ decisions will eventually be mapped into *routing* decisions from a network node to another.

It is our purpose in this paper to make these decisions *jointly*, accounting for the complex and often counterintuitive way they depend upon one another. To this end, we propose a queuing-based model, synthetically accounting for all the main aspects of the entities composing 5G networks. For physical hosts, we properly model their limited computational capabilities, and the fact that such capabilities must be split among the VNFs deployed at the same host. For VNFs, we account for the fact that they have minimum requirements in terms of computational capabilities, which have to match the

vertical requirements, and that, if additional capabilities are available, they will be exploited by VNFs to run faster. This leads to a flexible CPU allocation – an aspect that has not been addressed by previous work. Also, we tackle the case where, due to high traffic load, multiple instances of the same VNF are needed. As far as links between hosts are concerned, we take into account their finite capacity and non-zero delay. Finally, we model the fact that different classes of service requests, with different quality of service (QoS) requirements, coexist in the network and may share a portion of their VNF graphs.

We adopt a queuing model owing to the nature of 5G traffic. Indeed, a substantial fraction of this traffic, most notably coming from Internet-of-things (IoT) and machine-to-machine (M2M) applications, will consist of requests (customers) that originate from clients and then go through one or more computational stages (queues), triggering additional requests in the process. Queue networks are also the natural way to model the interaction pattern supported by present-day technologies such as Amazon Lambda and Amazon Greengrass [2], and endorsed by ETSI specifications (see the TASK request in [3, Sec. 6.11]).

Given our system model, we take *latency* as the main key performance indicator (KPI) and formulate an optimization problem that minimizes the ratio between the actual and the maximum allowed latency, across all service classes. Such a problem is impractical to solve directly, owing to its overwhelming complexity; thus, we present an efficient solution strategy, leveraging on *sequential* decision making. It is worth stressing that sequentially does not mean separately: decisions are made one after the other, but the interaction between them is still accounted for.

We evaluate the performance of our strategy against state-of-the-art alternatives, as well as against the global optimum. Owing to the diverse types of services (many still to be envisioned) that 5G networks will serve, we perform our evaluation for several different VNF graphs, ranging from simple chains to meshed graphs and akin to those corresponding to real-world services, most notably virtual EPC (vEPC).

In summary, our main contributions are as follows:

- we model all the main components of 5G services and network slices, from VNFs to physical hosts and links;
- we allow VNFs to be connected in arbitrarily graphs, as opposed to simpler chains or direct acyclic graphs (DAGs), and account for the need to deploy multiple instances of the same VNF;
- as a unique feature of our model, we allow *flexible* CPU allocation decisions, accounting for the fact that the

¹In this paper, we only refer to CPU allocation for simplicity; notice however that our model and methodology can take into account any kind of resource, e.g., disk or RAM space.

same VNF placement can correspond to multiple CPU allocation strategies;

- we study how such allocation decisions influence the system performance when subsets of VNFs are *shared* by multiple services, with different QoS limits;
- we devise an efficient and effective heuristic to make joint VNF placement and CPU allocation decisions, and find that it consistently performs very close to the optimum throughout different VNF graphs;
- we state and prove several properties that optimal CPU allocation decisions have under full-load conditions, and use such properties to further simplify the solution process.

The remainder of the paper is organized as follows. Sec. II reviews related work, highlighting the novelty of our contribution. Sec. III and Sec. IV introduce the system model and problem formulation, and analyze the complexity of the latter. Sec. V presents our solution concept, while Sec. VI describes how we deal with the special case of full-load conditions. Sec. VII addresses scenarios with multiple VNF instances. Finally, Sec. VIII presents performance evaluation results and Sec. IX concludes the paper.

II. RELATED WORK

Network-centric optimization. Many works, including [4]–[8], tackle the problems of VNF placement and routing from a network-centric viewpoint, i.e., they aim at minimizing the load of network resources. In particular, [4] seeks to balance the load on links and servers, while [5] studies how to optimize routing to minimize network utilization. The above approaches formulate mixed-integer linear programming (MILP) problems and propose heuristic strategies to solve them. [6], [7] and [8] formulate ILP problems, respectively aiming at minimizing the cost of used links and network nodes, minimizing resource utilization subject to QoS requirements, and minimizing bitrate variations through the VNF graph.

Service provider’s perspective. Several recent works take the viewpoint of a service provider, supporting multiple services that require different, yet overlapping, sets of VNFs, and seek to maximize its revenue. The works [9], [10] aim at minimizing the energy consumption resulting from VNF placement decisions. [11], [12] study how to place VNFs between network-based and cloud servers so as to minimize the cost, and [13] studies how to design the VNF graphs themselves, in order to adapt to the network topology.

User-centric perspective. Closer to our own approach, several works take a user-centric perspective, aiming at optimizing the user experience. [14], [15] study the VNF placement problem, accounting for the computational capabilities of hosts as well as network delays. In [16], the authors consider inter-cloud latencies and VNF response times, and solve the resulting ILP through an affinity-based heuristic.

Virtual EPC. The Evolved Packet Core (EPC) is a prime example of a service that can be provided through SDN/NFV. Interestingly, different works use different VNF graphs to

implement EPC, e.g., splitting user- and control-plane entities [17]–[19] or joining together the packet and service gateways (PGW and SGW) [20], [21]. Our model and algorithms work with any VNF graph, which allows us to model any real-world service, including all implementations of vEPC.

A. Novelty

The closest works to ours, in terms of approach and/or methodology, are [14], [15], [16], and [20].

However, [14], [15] and [20] model the assignment of VNFs to physical servers as a generalized assignment problem, a resource-constrained shortest path problem and a MILP problem, respectively. This implies that either a server has enough spare CPU capacity to offer a VNF, or it does not. Our queuing model, instead, is the first to account for the flexible allocation of CPU to the VNFs running on a host, e.g., the fact that VNFs will work faster if placed at a scarcely-utilized server. Furthermore, [14] and [20] have as objective the minimization of costs and server utilization, respectively. Our objective, instead, is to minimize the latency incurred by requests of different classes, which changes the solution strategy that can be adopted. Finally, [15] only considers VNF chains instead of generic graphs, and does not account for the possibility that the quantity of traffic changes across processing steps.

The queuing model used in [16] is similar (in principle) to ours; however, [16] does not address overlaps between VNF graphs and only considers DAGs. Furthermore, in both [15] and [16] no CPU allocation decisions are made, and the objective is to minimize a global metric, ignoring the different requirements of different service classes. Finally, the affinity-based placement heuristic proposed in [16] neglects the inter-host latencies and this, as confirmed by our numerical results in Sec. VIII, can yield suboptimal performance.

III. SYSTEM MODEL

We model VNFs as $M/M/1$ queues, $q \in \mathcal{Q}$, whose customers correspond to service requests. The *class* each customer belongs to corresponds to the service each request is associated to; we denote the set of such classes by \mathcal{K} . The *service rate* $\mu(q)$ of each queue q reflects the amount of CPU each VNF is assigned to, and therefore influences the time it takes to process one service request.

Arrival rates at queue $q \in \mathcal{Q}$ are denoted by $\lambda_k(q)$. Note that these values are class-specific, and reflect the amount of traffic of different services. Class-specific *transfer probabilities* $\mathbb{P}(q_2|q_1, k)$ indicate the probability that a customer of class k enters VNF q_2 after being served by VNF q_1 . We also indicate with $\mathbb{P}(q|o, k)$ the probability that a request of class k starts its processing at VNF q .

Physical *hosts* are represented by elements $h \in \mathcal{H}$. Each host h has a finite *CPU capacity* κ_h . Host-specific κ_h values account for both different capabilities and different hosts, and the fact that some hosts may be assigned a low-power CPU state [22] for energy-saving purposes. This implies that energy

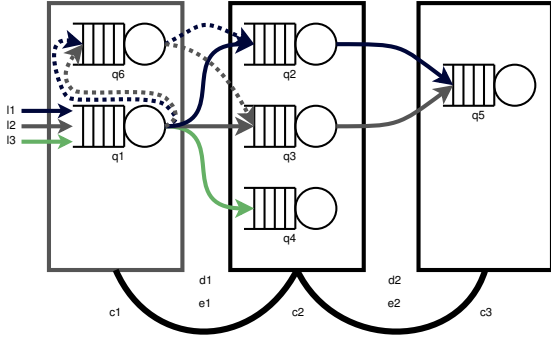


Fig. 1. Example 1: three service graphs, and six VNFs, corresponding to the six queues, placed across three physical hosts. Dashed and dotted lines represent the different paths that service requests can take.

constraints can be accounted for by properly setting the values of the κ_h parameters.

Going from host $h \in \mathcal{H}$ to host $l \in \mathcal{H}$ entails a deterministic *network latency* $\delta(h, l)$, which depends on the data transfer time between the two hosts. Furthermore, the link between hosts h and l has a finite capacity $C(h, l)$.

For clarity, an example on the use of the above notation is provided below.

Example 1: Assume that the network has to support three services: video streaming, gaming, and vehicle collision detection. Then the set of service classes is $\mathcal{K} = \{\text{video, game, veh}\}$. Each corresponds to a VNF graph, i.e.,

- *video streaming*: firewall – transcoder – billing;
- *gaming*: firewall – game server – billing;
- *vehicle collision detection*: firewall – collision detector.

Suspicious-looking packets belonging to the video streaming and gaming services can further be routed through a deep packet inspection (DPI) VNF. Hence, $\mathcal{Q} = \{\text{firewall, transcoder, billing, game server, collision detector, DPI}\}$.

There are three physical hosts $\mathcal{H} = \{h, l, m\}$, connected to each other through links characterized by a latency δ and a link capacity C . Fig. 1 illustrates the above quantities and shows a possible VNF placement across the three hosts. Routing can be deterministic, e.g., $\mathbb{P}(\text{billing}|\text{transcoder, video})=1$ or it can be probabilistic, e.g., $\mathbb{P}(\text{DPI}|\text{firewall, gaming})=0.1$ and $\mathbb{P}(\text{game server}|\text{firewall, gaming})=0.9$.

IV. PROBLEM FORMULATION AND COMPLEXITY

This section presents our joint problem formulation, including both VNF placement and CPU allocation decisions. For simplicity, we first formulate the problem for scenarios where exactly one instance of each VNF is to be deployed in the network; we will discuss the general case where multiple instances of the same VNF can be deployed in Sec. VII.

Decisions and decision variables. We have two main decision variables: a binary variable $A(h, q) \in \{0, 1\}$ represents whether VNF $q \in \mathcal{Q}$ is deployed at host $h \in \mathcal{H}$, and a real variable $\mu(q)$ expresses the amount of CPU assigned to

VNF $q \in \mathcal{Q}$. Notice how $\mu(q)$ maps onto the service rate of the corresponding queue.

System constraints. As mentioned above, we present our model in the case where there is exactly one instance of each VNF deployed in the system. This translates into imposing: $\sum_{h \in \mathcal{H}} A(h, q) = 1, \quad \forall q \in \mathcal{Q}$. Additionally, we have to honor the computational capacity limits of physical hosts, i.e.,:

$$\sum_{q \in \mathcal{Q}} A(h, q) \mu(q) \leq \kappa_h, \quad \forall h \in \mathcal{H}. \quad (1)$$

Arrival rates and system stability. Recall that input parameters $\lambda_k(q)$ express the rate at which *new* requests of service class k arrive at queue $q \in \mathcal{Q}$. We can then define an auxiliary variable $\hat{\lambda}_k(q)$, expressing the *total* rate of requests of class k that enter queue q , either from outside the system or from other queues. For any $k \in \mathcal{K}$, we have:

$$\hat{\lambda}_k(q) = \mathbb{P}(q|\circ, k) \sum_{q \in \mathcal{Q}} \lambda_{k,q} + \sum_{p \in \mathcal{Q} \setminus \{q\}} \mathbb{P}(q|p, k) \hat{\lambda}_k(p). \quad (2)$$

We can then define another auxiliary variable $\Lambda(q)$, expressing the total arrival rate of requests of any class entering queue q :

$$\Lambda(q) = \sum_{k \in \mathcal{K}} \hat{\lambda}_k(q).$$

Using $\Lambda(q)$, we can impose *system stability*, requesting that, for each queue, the arrival rate does not exceed the service rate:

$$\Lambda(q) < \mu(q), \quad \forall q \in \mathcal{Q}. \quad (3)$$

In other words, each VNF should receive *at least* enough CPU to deal with the incoming traffic. If additional CPU is available at the host, it will be exploited to further speed up the processing of requests.

Latency. This is our main metric of interest. It is due to two components: the processing latency and the network transfer latency.

The processing time, i.e., the time it takes for a service request of class k to traverse VNF q is represented by an auxiliary variable $R_k(q)$. For FCFS (first come, first serve) and PS (processor sharing) queuing disciplines, we have:

$$R_k(q) = \frac{1}{\mu(q) - \Lambda(q)}, \quad \forall q \in \mathcal{Q} \quad (4)$$

Note that the right-hand side of (4) does not depend on class k ; intuitively, this is because the queuing disciplines we consider are unaware of traffic classes. The response times for other queuing disciplines, including those accounting for priority levels and/or preemption, cannot be expressed in closed form. It is also worth stressing that present-day implementations of multi-access edge computing (MEC) [2] are based on FIFO discipline, and do not support preemption.

To compute the network latency that requests incur when transiting between hosts, we first need the expected number of times, $\gamma_k(q)$, that a request of class k visits VNF $q \in \mathcal{Q}$, i.e.,

$$\gamma_k(q) = \mathbb{P}(q|\circ, k) + \sum_{p \in \mathcal{Q} \setminus \{q\}} \mathbb{P}(q|p, k) \gamma_k(p). \quad (5)$$

In the right-hand side of (5), the first term is the probability that requests start their processing at queue q , and the second is the probability that requests arrive there from another queue p . Note that $\gamma_k(q)$ is not an auxiliary variable, but a quantity that can be computed offline given the transfer probabilities \mathbb{P} . Using $\gamma_k(q)$, the network latency incurred by requests of class k is:

$$\sum_{q,r \in \mathcal{Q}} \gamma_k(q) \mathbb{P}(r|q,k) \sum_{h,l \in \mathcal{H}} \delta(h,l) A(h,q) A(l,r). \quad (6)$$

We can read (6) from left to right, as follows. Given a service request of class k , it will be processed by VNF q for $\gamma_k(q)$ number of times. Every time, it will move to VNF r with probability $\mathbb{P}(r|q,k)$. So doing, it will incur latency $\delta(h,l)$ if q and r are deployed at hosts h and l , respectively (i.e., if $A(h,q) = 1$ and $A(l,r) = 1$).

The average total latency of requests of the generic class k is therefore given by:

$$D_k = \sum_{q \in \mathcal{Q}} R_k(q) + \sum_{q,r \in \mathcal{Q}} \gamma_k(q) \mathbb{P}(r|q,k) \sum_{h,l \in \mathcal{H}} A(h,q) A(l,r) \delta(h,l).$$

Link capacity. Given the finite link capacity $C(h,l)$, which limits the number of requests that move from any VNF at host h to any VNF at host l , we have:

$$\sum_{k \in \mathcal{K}} \sum_{q,r \in \mathcal{Q}} \hat{\lambda}_k(q) \mathbb{P}(r|q,k) A(q,h) A(r,l) \leq C(h,l). \quad (7)$$

Constraint (7) contains a summation over all classes k and all VNFs $q, r \in \mathcal{Q}$, such that q is deployed at h and r is deployed at l , as expressed by the A -variables. For each of such pair of VNFs, $\hat{\lambda}_k(q)$ is the rate of the requests of class k that arrive at q . Multiplying it by $\mathbb{P}(r|q,k)$, we get the rate at which requests move from VNF q to VNF r , hence from host h to host l .

Objective. D_k defined above represents the average latency incurred by requests of class k . In our objective function, we have to combine these values in a way that reflects the differences between such classes, most notably, their different QoS limits. Thus, we consider for each class k the *ratio* of the actual latency D_k to the limit latency D_k^{QoS} , and seek to minimize the maximum of such ratios:

$$\min_{A, \mu} \max_{k \in \mathcal{K}} \frac{D_k}{D_k^{\text{QoS}}}. \quad (8)$$

Importantly, the above objective function not only ensures fairness among service classes while accounting for their limit latency, but it also guarantees that the optimal solution will match *all* QoS limits if possible. More formally:

Property 1: If there is a non-empty set of solutions that meet constraints (1)–(7) and honor the services QoS limits, then the optimal solution to (8) falls in such a set.

The proof is omitted for brevity and can be found in [23].

Furthermore, when no solution meeting all QoS limits exists, the the solution optimizing (8) will minimize the damage by keeping all latencies as close as possible to their limit values.

A. Problem complexity

The problem of optimizing (8) subject to constraints (1)–(7) includes both binary ($A(h,q)$) and continuous ($\mu(q)$) variables. More importantly, constraint (7) and objective (8) are nonlinear and non-convex, as both include products between variables. The binary part of the problem is akin to the max-cut problem in graph theory [24], which has been proven to be NP-hard [25]. Indeed, our problem is even harder, as it includes evaluating and optimizing a non-convex function.

Such overwhelming complexity rules out not only the possibility to directly optimize the problem through a solver, but also commonplace solution strategies based on *relaxation*, i.e., allowing binary variables to take values anywhere in $[0, 1]$. Even if we relaxed the $A(h,q)$ variables, we would still be faced with a non-convex formulation, for which no algorithm is guaranteed to find a global optimum. We therefore present an efficient, *decoupled* solution strategy, leveraging on sequential decision making.

V. SOLUTION STRATEGY

Our solution strategy is based on decoupling the problems of VNF placement and CPU allocation, and then sequentially and iteratively making these decisions. We begin by presenting our VNF placement heuristic, called *MaxZ*, in Sec. V-A, and then discuss CPU allocation in Sec. V-B.

A. The MaxZ placement heuristic

As mentioned earlier, the two main sources of problem complexity are binary variables and non-convex functions in both objective (8) and constraint (7). Our VNF placement heuristic walks around these issues by:

- 1) formulating a convex version of the problem;
- 2) solving it through an off-the-shelf solver;
- 3) computing, for each VNF q and host h , a *score* $Z(h,q)$, expressing how confident we feel about placing q in h ;
- 4) considering the maximum score $Z(h^*, q^*)$ and placing VNF q^* at host h^* ;
- 5) repeating steps 2–4 until all VNFs are placed.

The name of the heuristic comes from step 4, where we seek for the highest score Z .

1) *Steps 1–2: convex formulation:* In order to make the problem formulation in Sec. IV convex, first we need to get rid of binary variables; specifically, we replace the binary variables $A(h,q) \in \{0, 1\}$ with continuous variables $\tilde{A}(h,q) \in [0, 1]$.

We also need to remove the products between \tilde{A} -variables (e.g., in (6) and (7)), by replacing them with a new variable. To this end, for each pair of VNFs q and r and physical hosts h and l , we introduce a new variable $\Phi(h,l,q,r) \in [0, 1]$, and impose that:

$$\Phi(h,l,q,r) \leq \tilde{A}(h,q), \quad \forall h,l \in \mathcal{H}, q,r \in \mathcal{Q}; \quad (9)$$

$$\Phi(h,l,q,r) \leq \tilde{A}(l,r), \quad \forall h,l \in \mathcal{H}, q,r \in \mathcal{Q}; \quad (10)$$

$$\Phi(h,l,q,r) \geq \tilde{A}(h,q) + \tilde{A}(l,r) - 1, \quad \forall h,l \in \mathcal{H}, q,r \in \mathcal{Q}. \quad (11)$$

The intuition behind constraints (9)–(11) is that $\Phi(h, l, q, r)$ mimics the behavior of the product $\tilde{A}(h, q)\tilde{A}(l, r)$: if either $\tilde{A}(h, q)$ or $\tilde{A}(l, r)$ are close to 0, then (9) and (10) guarantee that $\Phi(h, l, q, r)$ will also be close to zero; if both values are close to one, then (11) allows also $\Phi(h, l, q, r)$ to be close to one.

Another product between variables, i.e., a term in the form $A(h, q)\mu(q)$, appears in (1). Following a similar approach, we introduce a set of new variables, $\psi(h, q)$, expressing the fraction of h 's CPU that is allocated to q . We then impose:

$$\psi(h, q) \leq \tilde{A}(h, q), \quad \forall h \in \mathcal{H}, q \in \mathcal{Q}; \quad (12)$$

$$\sum_{q \in \mathcal{Q}} \psi(h, q) \leq 1, \quad \forall h \in \mathcal{H}, \quad (13)$$

which mimic (1). By replacing all products between \tilde{A} -variables with a Φ -variable and all products between \tilde{A} - and μ -variables with a ψ -variable, we obtain a *convex problem*, which can efficiently be solved through commercial solvers.

2) *Steps 3–4: Z-score and decisions*: Let us assume that no VNF has been placed yet. We then solve an instance of the convex problem described in Sec. V-A1, and use the values of the variables $\tilde{A}(h, q)$ and $\psi(h, q)$ to decide which VNF to place at which host.

Recall that $\tilde{A}(h, q)$ is the relaxed version of our placement variable $A(h, q)$, so we would be inclined to use that to make our decision. However, we also need to account for how much computational capacity VNFs would get, as expressed by $\psi(h, q)$. If such a value falls below the threshold $T_\psi(h, q) = \frac{\Lambda(q)}{\kappa_h}$, then VNF q may not be able to process the incoming requests, i.e., constraint (3) may be violated.

To prevent this, we define our Z-score, i.e., how confident we are about placing VNF q at host h , as follows:

$$Z(h, q) = \tilde{A}(h, q) + \mathbf{1}_{\psi(h, q) \geq T_\psi(h, q)}, \quad (14)$$

where $\mathbf{1}$ is the indicator function. Recalling that \tilde{A} -values are constrained between 0 and 1, favoring high values of (14) means that we prefer a deployment that results in ψ -values greater than the threshold, if such a deployment exists. Otherwise, we make the placement decision based on the \tilde{A} -values only.

Specifically, we select the host h^* and VNF q^* associated with the maximum Z, i.e., $h^*, q^* \leftarrow \arg \max_{h \in \mathcal{H}, q \in \mathcal{Q}} Z(h, q)$, and place VNF q^* in host h^* . We fix this decision and repeat the procedure till all VNFs are placed (i.e., we perform exactly $|\mathcal{Q}|$ iterations).

We now present two example runs of MaxZ, for two scenarios with different inter-host latencies.

Example 2: Consider a simple case with two hosts $\mathcal{H} = \{h_1, h_2\}$ with the same CPU capacity $\kappa_h = 10$ requests/s, two VNFs $\mathcal{Q} = \{q_1, q_2\}$, and only one request class k with $\lambda_k = 1$ requests/s. Requests need to subsequently traverse q_1 and q_2 . The inter-host latency $\delta(h_1, h_2)$ is set to

10 ms, while $D^{\text{QoS}} = 100$ ms. Then, intuitively, the optimal solution is to deploy one VNF per host.

We solve the problem in Sec. V-A1. After the first iteration, we obtain $\tilde{\mathbf{A}} = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$, $\psi = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$, and $\mathbf{Z} = \begin{bmatrix} 1.5 & 1.5 \\ 1.5 & 1.5 \end{bmatrix}$ ². In such a case, using a tie-breaking rule, we place VNF q_1 at host h_1 . In the second iteration, we have $\tilde{\mathbf{A}} = \begin{bmatrix} 1 & 0.38 \\ 0 & 0.62 \end{bmatrix}$, $\psi = \begin{bmatrix} 0.8 & 0.19 \\ 0 & 0.61 \end{bmatrix}$, and $\mathbf{Z} = \begin{bmatrix} 2 & 1.38 \\ 0 & 1.62 \end{bmatrix}$. We ignore the entries pertaining to VNF q_1 that has been already placed and, since $Z(h_2, q_2) > Z(h_1, q_2)$, we deploy VNF q_2 at host h_2 , which corresponds to the intuition that, given the small value of δ , VNFs should be spread across the hosts.

Example 3: Let us now consider the same scenario as in Example 2, but assume a much longer latency $\delta(h_1, h_2) = 200$ ms. The best solution will now be to place both VNFs at the same host.

After the first iteration, we obtain $\tilde{\mathbf{A}} = \begin{bmatrix} 0.7 & 0.7 \\ 0.3 & 0.3 \end{bmatrix}$, $\psi = \begin{bmatrix} 0.5 & 0.5 \\ 0.3 & 0.3 \end{bmatrix}$, and $\mathbf{Z} = \begin{bmatrix} 1.7 & 1.7 \\ 1.3 & 1.3 \end{bmatrix}$. Again using a tie-breaking rule, we place VNF q_1 at host h_1 . In the second iteration, we have $\tilde{\mathbf{A}} = \begin{bmatrix} 1 & 0.2 \\ 0 & 0.2 \end{bmatrix}$, $\psi = \begin{bmatrix} 0.6 & 0.4 \\ 0 & 0.2 \end{bmatrix}$, and $\mathbf{Z} = \begin{bmatrix} 2 & 1.8 \\ 0 & 1.2 \end{bmatrix}$. We again ignore the entries in the first column and, since $Z(h_1, q_2) > Z(h_2, q_2)$, we place VNF q_2 at host h_1 , making optimal decisions.

B. CPU allocation

Once the MaxZ heuristic introduced in Sec. V-A provides us with deployment decisions, we need to decide the CPU allocation, i.e., the values of the $\mu(q)$ variables in the original problem described in Sec. IV. This can be achieved simply by solving the problem in (8) but keeping the deployment decision fixed, i.e., replacing the $A(h, q)$ variables with parameters whose values come from the MaxZ heuristic. Indeed, we can prove the following property.

Property 2: If the deployment decisions are fixed, then the problem of optimizing (8) subject to (1)–(7) is convex.

The proof can be found in our technical report [23].

Property 2 guarantees that we can make our CPU allocation decisions, i.e., decide on the $\mu(q)$ values, in polynomial time. We can further enhance the solution efficiency by reducing the optimization problem to the resolution of a system of equations, through the Karush-Kuhn-Tucker (KKT) conditions.

1) *KKT conditions*: In order to derive the KKT conditions for the problem stated in Sec. IV, we need to re-write objective (8) in standard form. This requires introducing an auxiliary variable ρ representing the maximum $\frac{D_k}{D_k^{\text{QoS}}}$ ratio, and imposing that for each class $k \in \mathcal{K}$:

$$\rho \geq \frac{1}{D_k^{\text{QoS}}} \left(\sum_{q \in \mathcal{Q}} \frac{\gamma_k(q)}{\mu(q) - \Lambda(q)} + \sum_{h_1, h_2 \in \mathcal{H}} \nu_k(h_1, h_2) \delta(h_1, h_2) \right) \quad (15)$$

where $\nu_k(h_1, h_2)$ is the expected number of times that a request of class k ever travels from host h_1 to host h_2 . These quantities depend upon the placement decisions $A(h, q)$ and

²In all matrices, rows correspond to hosts and columns to VNFs.

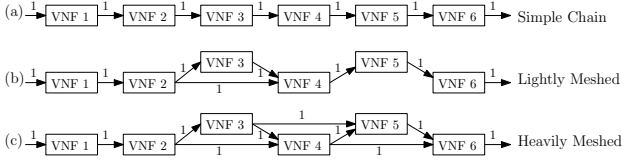


Fig. 2. The VNF graphs we consider in our performance evaluation, reflecting real-world service implementations.

are input parameters to the CPU allocation problem. At this point, the objective is simply to minimize ρ .

We also re-write constraints (1), (3) and (15) in normal form, and associate to them the multipliers M_q , M_h and M_k respectively. The resulting Lagrangian function is:

$$L = \rho + \sum_{q \in \mathcal{Q}} M_q X_q + \sum_{h \in \mathcal{H}} M_h Y_h + \sum_{k \in \mathcal{K}} M_k W_k, \quad (16)$$

where:

$$X_q = -\mu(q) + \Lambda(q);$$

$$Y_h = \sum_{q \in \mathcal{Q}} A(h, q) \mu(q) - \kappa_h;$$

$$W_k = \sum_{q \in \mathcal{Q}} \frac{\gamma_k(q)}{D_k^{\text{QoS}}} \frac{1}{\mu(q) - \Lambda(q)} + \sum_{h_1, h_2 \in \mathcal{H}} \nu_k(h_1, h_2) \frac{\delta(h_1, h_2)}{D_k^{\text{QoS}}} - \rho.$$

The first necessary KKT condition, i.e., $\nabla_{r\rho\mu(q)} L = 0$, translates into the following equations:

$$\frac{\partial}{\partial \rho} L = 0 \iff 1 - \sum_{k \in \mathcal{K}} M_k = 0. \quad (17)$$

Furthermore, for each $q \in \mathcal{Q}$, we have:

$$\begin{aligned} \frac{\partial}{\partial \mu(q)} L = 0 \iff & -M_q + \sum_{h \in \mathcal{H}} M_h A(h, q) + \\ & - \sum_{k \in \mathcal{K}} \frac{M_k \gamma_k(q)}{D_k^{\text{QoS}}} \frac{1}{(\mu(q) - \Lambda(q))^2} = 0 \end{aligned} \quad (18)$$

Finally, complementary slackness requires that either the inequality constraints are active, or the corresponding multipliers are zero, i.e.,

$$M_q X_q = 0, \quad \forall q \in \mathcal{Q}, \quad (19)$$

$$M_h Y_h = 0, \quad \forall h \in \mathcal{H}, \quad (20)$$

$$M_k W_k = 0, \quad \forall k \in \mathcal{K}. \quad (21)$$

Based on (19)–(21), the multipliers assume the following meaning:

- M_q is zero for all *stable* queues, i.e., the queues fulfilling the condition $\mu(q) > \Lambda(q)$;
- M_h is zero for all *non-strained* hosts, i.e., hosts used for less than their CPU capacity κ_h ;
- M_k is zero for all *non-critical* classes, i.e., classes for which the $\frac{D_k}{D_k^{\text{QoS}}}$ ratio is strictly lower than ρ .

We can now determine the *global* computational complexity of our approach, including the VNF placement through the MaxZ heuristic and the CPU allocation by optimizing (8).

Property 3: Our solution strategy, including the MaxZ VNF placement heuristic in Sec. V-A and the CPU allocation strategy in Sec. V-B has polynomial computational complexity. The proof can be found in [23].

VI. SPECIAL CASE: FULL-LOAD CONDITIONS

In this section, we seek to further reduce the complexity of the CPU allocation problem. Let us start from the Lagrange multipliers derived earlier, and recall that we require *stability*, i.e., $\Lambda(q) < \mu(q)$, hence $M_q = 0$ for all queues $q \in \mathcal{Q}$;

Given the above and (18), we can write that, for each queue $q \in \mathcal{Q}$ deployed at host h ,

$$M_h = \sum_{k \in \mathcal{K}} M_k \frac{\gamma_k(q)}{D_k^{\text{max}}} \frac{1}{(\mu(q) - \Lambda(q))^2}. \quad (22)$$

Recalling the meaning of the multipliers, we can state the following lemma and properties. Proofs are omitted for brevity and reported in [23].

Lemma 1: If CPU assignment decisions are made optimizing the objective (8), then there exists at least one critical class, i.e., for which equality holds in (15).

Property 4: All hosts traversed by service requests of critical classes are strained.

Property 5: VNFs deployed at a strained host serve at least one critical class each.

In summary, there is at least one critical class, *all* the hosts it traverses are strained, and *each* of the VNFs deployed on the strained hosts (not only the ones serving requests of the original critical class) serve at least one critical class. This can lead to a cascade effect, as shown in Example 4.

Example 4: Consider the case in Fig. 1. By Lemma 1, at least one class is critical; let us assume that such a class is collision detection. From Property 4, all hosts traversed by collision detection requests, i.e., hosts h and l , are strained. Since host l is strained, from Property 5 it follows that each of its VNFs serves at least one critical class. Since the transcoder queue only serves the video class, the video class is critical. Similarly, since the game server only serves the game class, that class is critical as well. Finally, both video and game classes traverse host m ; therefore, by Property 4, that host is critical as well.

In scenarios like the one in Example 4, where all classes are critical and all hosts are strained, we have:

$$\sum_{q \in \mathcal{Q}} \frac{\gamma_k(q)}{D_k^{\text{QoS}}} \frac{1}{\mu(q) - \Lambda(q)} + \sum_{h_1, h_2 \in \mathcal{H}} \nu_k(h_1, h_2) \frac{\delta(h_1, h_2)}{D_k^{\text{QoS}}} = \rho \quad \forall k \in \mathcal{K},$$

$$\sum_{q \in \mathcal{Q}} A(h, q) \mu(q) = \kappa_h, \quad \forall h \in \mathcal{H}.$$

The above equalities can be combined with the KKT conditions stated in Sec. V-B1, thus forcing $Y_h = 0 \quad \forall h \in \mathcal{H}$ and $W_k = 0 \quad \forall k \in \mathcal{K}$. This greatly simplifies and speeds up the

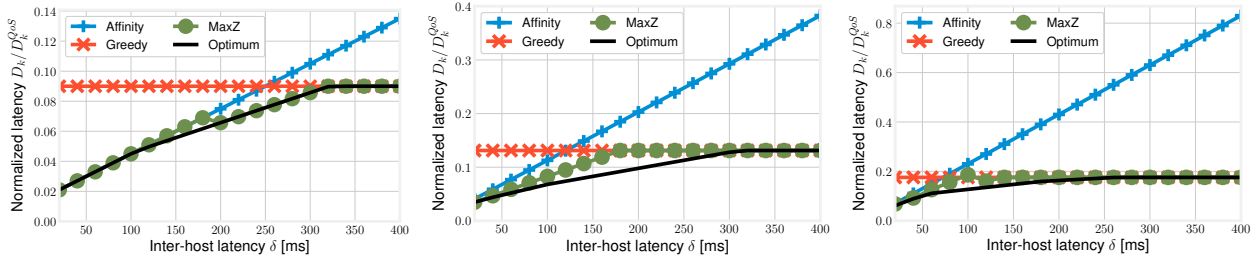


Fig. 3. Normalized latency as a function of inter-host latency δ for the chain (left), light mesh (center), heavy mesh (right) VNF graphs. Note that the y-axis scale varies across the plots.

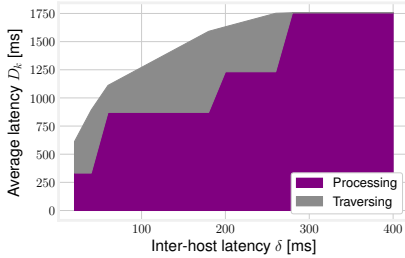


Fig. 4. Breakdown of the total latency as a function of inter-host latency δ , for the heavy mesh topology and the MaxZ deployment strategy.

process of finding the optimal CPU allocation values $\mu(q)$. Even more importantly, the equalities above also simplify the computation of the $\mu(q)$ values within each iteration of the MaxZ heuristic described in Sec. V-A.

VII. MULTIPLE VNF INSTANCES

So far, we presented our system model and solution strategy in the case where exactly one instance of each VNF has to be deployed. This is not true in general; some VNFs may need to be replicated owing to their complexity and/or load.

If the number N_q of instances of VNF q to be deployed is known, then we can replace VNF q in the VNF graph with N_q replicas thereof, labeled q^1, q^2, \dots, q^{N_q} , each with the same incoming and outgoing edges. With regard to the $\Lambda(q)$ requests/s that have to be processed by *any* instance of VNF q , they are split among the instances. If $f(q, i)$ is the fraction of requests for VNF q that is processed by instance q^i (and thus $\sum_{i=1}^{N_q} f(q, i) = 1$), then instance q^i gets requests at a rate $f(q, i)\Lambda(q)$. It is important to stress that once the $f(q, i)$ splitting fractions are known, then the resulting problem can be solved with the approach described in Sec. V.

Establishing the $f(q, i)$ values is a complex problem; indeed, straightforward solutions like uniformly splitting flows (i.e., $f(q, i) = \frac{1}{N_q}$), are in general suboptimal. We therefore resort to a *pattern search* [26] iterative approach.

Without loss of generality, we describe our approach in the simple case $N_q = 2$. In this case, the splitting values are $f(q, 1) = f$ and $f(q, 2) = 1 - f$. Given an initial guess f_0 , an initial step Δ and a minimum step $\epsilon < \Delta$, we proceed as follows:

- 1) we initialize the splitting factor f to the initial guess f_0 ;

- 2) using the procedure detailed in Sec. V, we compute the objective value (8) for the splitting values f , $f + \Delta$, and $f - \Delta$;
- 3) if the best result (i.e., the lowest value of (8)) is obtained for splitting value $f + \Delta$ or $f - \Delta$, then we replace f with that value and loop to step 2;
- 4) otherwise, we reduce Δ by half;
- 5) if now Δ is lower than ϵ , the algorithm terminates;
- 6) otherwise, we loop to step 2.

The intuition of the pattern-search procedure is similar, in principle, to gradient-search methods. If we find that using $f + \Delta$ or $f - \Delta$ instead of f produces a lower delay, then we replace the current value of f with the new one; otherwise, we try new f -values closer to the current one. When we are satisfied that there are no better f -values further than ϵ from the current one, the search ends.

Notice that in step 2 of our procedure we run the decision-making procedure described in Sec. V; this implies that, once we find the best value of f , we also know the best VNF placement and CPU allocation decisions.

VIII. NUMERICAL RESULTS

Reference scenario. We consider a reference topology with three hosts having CPU capacity $\kappa_h = 10$ requests/ms each. The hosts are connected by links having latency δ that varies between 50 ms and 400 ms. For simplicity, we disregard the link capacity, i.e., we assume computation to be the bottleneck in our scenario. Throughout our performance evaluation, we benchmark the MaxZ placement heuristic in Sec. V-A against the following alternatives:

- *global optimum*, found through brute-force search of all possible deployments;
- *greedy*, where the number of used hosts is minimized, i.e., VNFs are concentrated as much as possible;
- *affinity-based* [16], trying to place at the same host VNFs with high transition probability between them.

After the VNF placement decisions are made, we compute the optimal CPU allocation, i.e., the optimal $\mu(q)$ values, as explained in Sec. V-B. It is important to remark that the CPU allocation strategy is the same for all placement strategies.

We first focus on a single request class k , fix the arrival rate to $\lambda_k = 1$ requests/ms, and compare the three VNF graphs depicted in Fig. 2, ranging from a simple chain to a

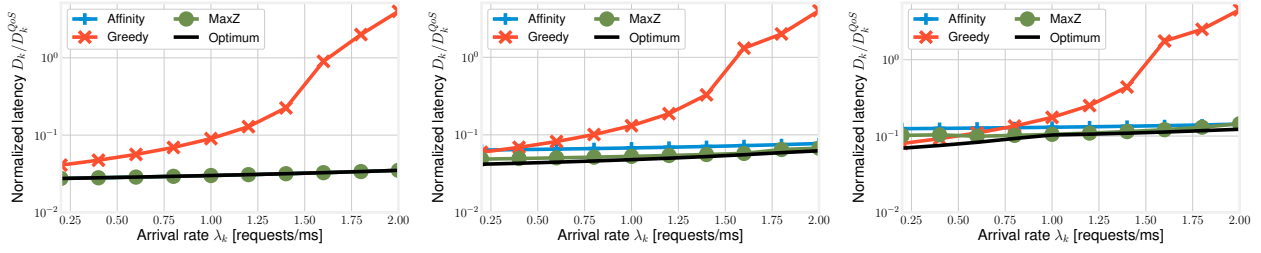


Fig. 5. Normalized latency (log scale) as a function of arrival rate λ for the chain (left), light mesh (center), heavy mesh (right) VNF graphs.

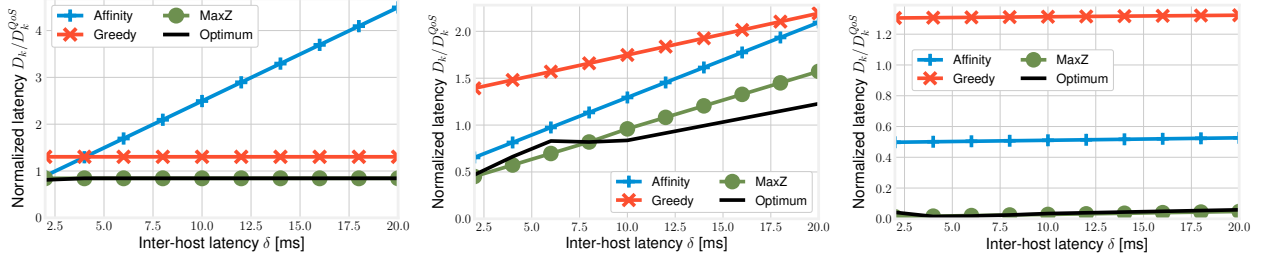


Fig. 6. Multi-class scenario, heavy mesh graph: normalized latency vs. arrival rate λ for the low-latency (left), medium-latency (center), high-latency (right) service classes. Note that the y-axis scale varies across the plots.

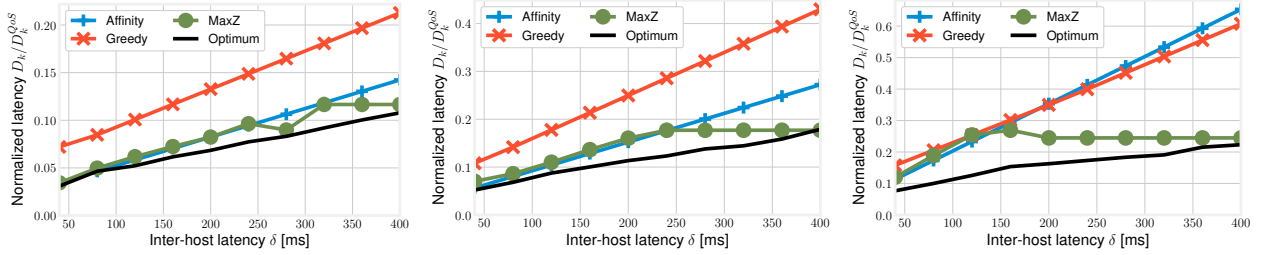


Fig. 7. Multi-instance scenario: normalized latency vs. inter-host latency δ for the chain (left), light mesh (center), heavy mesh (right) VNF graphs. Note that the y-axis scale varies across the plots.

complex meshed topology. Notice how in graphs (b) and (c) requests can *branch* and *merge*, i.e., the number of requests outgoing from a VNF does not match the number of incoming ones. This is the case with several real-world functions; in particular, the light and heavy mesh topologies are akin to vEPC implementations where user- and control-plane entities are joint [21] and split [18], [19].

Effect of the inter-host latency. Fig. 3 shows the average global latency as a function of the inter-host latency δ , for the VNF graphs presented in Fig. 2. We can observe that the performance of Greedy is always the same regardless of δ , as all VNFs are deployed at the same host. On the other hand, the performance of Affinity-based is quite good for low values of δ , but then quickly degrades, due to the fact that the affinity-based heuristic disregards inter-host latency. As far as MaxZ is concerned, it exhibits an excellent performance: it consistently yields a substantially lower latency compared to Greedy and Affinity-based, and is always very close to the optimum.

Fig. 4 focuses on the heavy mesh topology, and breaks down the total latency yielded by MaxZ into its computation

and traversing latency components. Processing latency only depends upon the VNF placement, while traversing latency depends upon both the VNF placement (which determines how many inter-host links are traversed) and the per-link latency δ . When δ is low, MaxZ tends to spread the VNFs across all available hosts, in order to assign more CPU. As δ increases, the placement becomes more and more concentrated (thus resulting in lower $\mu(q)$ values and higher processing times), until, when δ is very high, all VNFs are placed at the same host and there is no traversing latency at all.

Fig. 3 and Fig. 4 clearly illustrate the importance of flexible CPU allocation. If we only accounted for the minimum CPU required by VNFs, as in [14], [16], we could place all of them in the same host, as the Greedy strategy does. This would result, as we can see from the far right in Fig. 4, in high processing times and *two unused* hosts.

Effect of arrival rate. We now fix the inter-host latency to $\delta = 50$ ms, and change the arrival rate λ between 0.1 requests/ms and 2 requests/ms; Fig. 5 summarizes the latency yielded by the placement strategies we study.

A first observation concerns the Greedy strategy: since all VNFs are placed in the same host, as λ increases, VNFs receive an amount of CPU that is barely above the minimum limit $\Lambda(q)$. This, as per (4), results in processing times that grow very large. The difference between the other placement strategies tends to become less significant; intuitively, this is because processing times dominate the total latency, and thus spreading the VNFs as much as possible is always a sensible solution. MaxZ still consistently outperforms Affinity-based, and performs very close to the optimum.

Multi-class scenario. In Fig. 6, we move to a multi-class scenario where $|\mathcal{K}| = 3$ service classes share the same VNF graph. The three classes have limit latencies D^{QoS} of 10 ms, 45 ms, and 2 s, respectively corresponding to safety applications (e.g., collision detection), real-time applications (e.g., gaming), and delay-tolerant applications (e.g., video streaming). Fig. 6 shows that all placement strategies result in latencies that are roughly proportional to the limit ones. Also, the relative performance of the placement strategies remains unmodified – MaxZ outperforms Affinity-based and is close to the optimum, while Greedy yields much higher latency. Notice that, for very high values of δ , it is impossible to meet all QoS constraints, i.e., $\frac{D_k}{D_k^{\text{QoS}}} > 1$ for at least one class k . In these cases, MaxZ limits the damage by keeping the $\frac{D_k}{D_k^{\text{QoS}}}$ ratios as low as possible. It is also interesting to notice, in Fig. 6(center), that the latency yielded by MaxZ is actually lower than the optimum. This does not mean that MaxZ outperforms the optimum, rather that it gives a lower latency to the medium-latency, whose latency is closest to the QoS limit.

Multiple VNF instances. In Fig. 7, we drop the assumption that there is only one instance of each VNF; specifically, for VNF₄ and VNF₆ we allow two instances each. We can immediately see, by comparing Fig. 7 to Fig. 3, that allowing multiple VNF instances substantially decreases the total latency. More interestingly, we can observe that MaxZ always outperforms its alternatives, and is very close to the optimum, except for some cases when the topology is very complex.

IX. CONCLUSION

We presented a model for SDN/NFV-based 5G networks that is able to capture all their main features, including arbitrary VNF graphs, flexible CPU allocation to VNFs, and the possibility to have multiple instances of the same VNF. Leveraging this model, we presented a methodology, based on the MaxZ placement heuristic, to make joint VNF placement and CPU assignment decisions.

We combined MaxZ with a methodology to make optimal CPU allocation decisions, requiring to solve a convex optimization problem in the general case and a simple system of equations in full-load conditions. By evaluating our solution over several VNF graphs of different complexity, we consistently found it to outperform state-of-the-art alternatives and closely track optimal performance.

One direction for future work is enhancing the performance of our heuristic in multi-instance scenarios, by improving the

pattern search approach we adopted and further customizing it to our needs.

ACKNOWLEDGEMENT

This work is partially supported by the European Commission through the H2020 5G-TRANSFORMER project (Project ID 761536).

REFERENCES

- [1] NGMN Alliance, “Description of network slicing concept,” 2016.
- [2] Amazon. AWS Greengrass. <https://aws.amazon.com/greengrass/>.
- [3] ETSI. GS MEC 009: Mobile Edge Computing (MEC); General principles for Mobile Edge Service APIs.
- [4] A. Hirwe and K. Kataoka, “LightChain: A lightweight optimization of VNF placement for service chaining in NFV,” in *IEEE NetSoft*, 2016.
- [5] T. W. Kuo, B. H. Liou, K. C. J. Lin, and M. J. Tsai, “Deploying chains of virtual network functions: On the relation between link and server usage,” in *IEEE INFOCOM*, 2016.
- [6] A. Baumgartner, V. S. Reddy, and T. Bauschert, “Mobile core network virtualization: A model for combined virtual core network function placement and topology optimization,” in *IEEE NetSoft*, 2015.
- [7] F. Ben Jemaa, G. Pujolle, and M. Pariente, “Analytical Models for QoS-driven VNF Placement and Provisioning in Wireless Carrier Cloud,” in *ACM MSWiM*, 2016.
- [8] B. Addis, D. Belabed, M. Bouet, and S. Secci, “Virtual network functions placement and routing optimization,” in *IEEE CloudNet*, 2015.
- [9] A. Marotta and A. Kessler, “A Power Efficient and Robust Virtual Network Functions Placement Problem,” in *IEEE ITC*, 2016.
- [10] N. E. Khoury, S. Ayoubi, and C. Assi, “Energy-Aware Placement and Scheduling of Network Traffic Flows with Deadlines on Virtual Network Functions,” in *IEEE CloudNet*, 2016.
- [11] M. Mechtri, C. Ghribi, and D. Zeghlache, “A scalable algorithm for the placement of service function chains,” *IEEE Transactions on Network and Service Management*, 2016.
- [12] L. Gu, S. Tao, D. Zeng, and H. Jin, “Communication cost efficient virtualized network function placement for big data processing,” in *IEEE INFOCOM Workshops*, 2016.
- [13] J. Cao, Y. Zhang, W. An, X. Chen, J. Sun, and Y. Han, “VNF-FG design and VNF placement for 5G mobile networks,” *Science China Information Sciences*, 2017.
- [14] R. Cohen, L. Lewin-Eytan, J. S. Naor, and D. Raz, “Near optimal placement of virtual network functions,” in *IEEE INFOCOM*, 2015.
- [15] B. Martini, F. Paganelli, P. Capanera, S. Turchi, and P. Castoldi, “Latency-aware composition of virtual functions in 5G,” in *IEEE NetSoft*, 2015.
- [16] D. Bhamare, M. Samaka, A. Erbad, R. Jain, L. Gupta, and H. A. Chan, “Optimal virtual network function placement in multi-cloud service function chaining architecture,” *Computer Communications*, 2017.
- [17] A. Baumgartner, V. S. Reddy, and T. Bauschert, “Mobile core network virtualization: A model for combined virtual core network function placement and topology optimization,” in *IEEE NetSoft*, 2015.
- [18] G. Hasegawa and M. Murata, “Joint Bearer Aggregation and Control-Data Plane Separation in LTE EPC for Increasing M2M Communication Capacity,” in *IEEE GLOBECOM*, 2015.
- [19] A. Ksentini, M. Bagaa, and T. Taleb, “On Using SDN in 5G: The Controller Placement Problem,” in *IEEE Globecom*, 2016.
- [20] D. Dietrich, C. Papagianni, P. Papadimitriou, and J. S. Baras, “Network function placement on virtualized cellular cores,” in *COMSNETS*, 2017.
- [21] J. Prados-Garzon, J. J. Ramos-Munoz, P. Ameigeiras, P. Andres-Maldonado, and J. M. Lopez-Soler, “Modeling and Dimensioning of a Virtualized MME for 5G Mobile Networks,” *IEEE Transactions on Vehicular Technology*, 2017.
- [22] Intel. Power Management States: P-States, C-States, and Package C-States. <https://software.intel.com/en-us/articles/power-management-states-p-states-c-states-and-package-c-states>.
- [23] Proofs. <https://1drv.ms/b/s!A1VUQf67dIbfcSZH2bZfrOTKLQc>.
- [24] A. V. Goldberg and R. E. Tarjan, “A new approach to the maximum-flow problem,” *Journal of the ACM*, 1988.
- [25] C. H. Papadimitriou and M. Yannakakis, “Optimization, approximation, and complexity classes,” *Journal of computer and system sciences*, 1991.
- [26] R. M. Lewis and V. Torczon, “Pattern search methods for linearly constrained minimization,” *SIAM Journal on Optimization*, 2000.