

# Lyapunov-Guided Optimal Service Placement in Vehicular Edge Computing

Chaogang Tang<sup>1</sup>, Yubin Zhao<sup>2</sup>, Huaming Wu<sup>3,\*</sup>

<sup>1</sup> School of Computer Science and Technology, China University of Mining and Technology, Xuzhou 221116, China

<sup>2</sup> School of Microelectronics Science and Technology, Sun Yat-Sen University, Zhuhai 519082, China

<sup>3</sup> Center for Applied Mathematics, Tianjin University, Tianjin 300072, China

\* The corresponding author, email: whming@tju.edu.cn

**Cite as:** C. Tang, Y. Zhao, *et al.*, “Lyapunov-guided optimal service placement in vehicular edge computing,” *China Communications*, vol. 20, no. 3, pp. 201-217, 2023. DOI: 10.23919/JCC.2023.03.015

**Abstract:** Vehicular Edge Computing (VEC) brings the computational resources in close proximity to the service requestors and thus supports explosive computing demands from smart vehicles. However, the limited computing capability of VEC cannot simultaneously respond to large amounts of offloading requests, thus restricting the performance of VEC system. Besides, a mass of traffic data can incur tremendous pressure on the front-haul links between vehicles and the edge server. To strengthen the performance of VEC, in this paper we propose to place services beforehand at the edge server, e.g., by deploying the services/tasks-oriented data (e.g., related libraries and databases) in advance at the network edge, instead of downloading them from the remote data center or offloading them from vehicles during the runtime. In this paper, we formulate the service placement problem in VEC to minimize the average response latency for all requested services along the slotted timeline. Specifically, the time slot spanned optimization problem is converted into per-slot optimization problems based on the Lyapunov optimization. Then a greedy heuristic is introduced to the drift-plus-penalty-based algorithm for seeking the approximate solution. The simulation results reveal its advantages over others in

terms of optimal values and our strategy can satisfy the long-term energy constraint.

**Keywords:** vehicular edge computing; service placement; response latency; computational resources

## I. INTRODUCTION

Various advanced Information and Communication Technologies (ICT) are applied and integrated into Intelligent Transportation Systems (ITS) to pursue the harmonious unification of humans, vehicles and roads. Many benefits can be obtained from ITS, including guaranteeing traffic safety, raising transportation efficiency, and improving the transportation environment and driving experience. In this context, a wide variety of data can be gathered by mounted facilities in smart vehicles (e.g., SIM cards, GPS, sensors, and cameras). Furthermore, the consequent vehicular tasks have sprung up, which require computational resources and services for successful execution. Smart vehicles are usually equipped with computing facilities such as On-Board Unit (OBU) to satisfy the computational demands of these tasks. However, such “computers on wheels” can only offer limited computing capabilities, due to restricted physical size and energy supply.

To alleviate the contradiction between the skyrocketing demands for computational resources and the limited computing capabilities of vehicles themselves,

Received: Jan. 12, 2022

Revised: Mar. 15, 2022

Editor: Lin Gao

---

it becomes increasingly dominant that computational resources and services are provisioned outside the vehicles. For example, Vehicular Cloud Computing (VCC) and Vehicular Edge Computing (VEC) are two newly emergent computing paradigms in recent years, which are intended for addressing the above contradiction. The difference between VCC and VEC lies in that VEC shifts the computational capability from the remote data center to the network edge and provides computing services in close proximity to the resource requestors (i.e., vehicles). Comparatively speaking, VEC can better cater to the delay-strict requirements of vehicle tasks. In particular, vehicles offload their tasks/applications to the edge server for execution in VEC and the edge server is usually deployed at the Road-Side Unit (RSU) in a densely populated area. However, due to explosive growth in vehicular tasks and applications, tremendous pressure has been imposed on the computing capabilities of the edge server as well as the front-haul links between vehicles and an accessible computational access point, e.g., RSU. This is because the amount of computational resources in VEC is not unlimited compared to that at the cloud center.

Accordingly, we propose to place services in advance at the edge server. Service placement, also termed service deployment, strives to deploy the services/tasks-oriented data (e.g., related libraries and databases) in advance at the network edge (e.g., RSUs), instead of downloading them from a remote data center during the runtime or rush hours. When a vehicle sends its offloading request to the edge server, the execution result can be directly retrieved, if the corresponding service has been placed at the edge server. Otherwise, the corresponding service may be offloaded from the vehicle or downloaded from the data center, which increases both the transmission and calculation delay. As a result, the response delay is much longer, compared to the case when the service is cached.

Although there is extensive literature on task offloading and service placement in VEC [1–4], they seldom consider factors that may increase the difficulty of service placement. Such factors usually include limited sojourn time for vehicles, limited wireless coverage and vehicular offloading requests featured by temporal and spatial variation. In addition, most of the current works on service caching assume that service

caching does not incur energy consumption at the edge server, which, however, does not always hold in the task-oriented caching scenarios [5, 6]. In contrast, this paper does not only take into account the above factors but also the evaluation of service placement in the long run.

In this paper, we strive to make service placement decisions with a time effect and strategic update in a scenario where vehicles are covered by multiple RSUs owing to an increasing density of RSUs. Specifically, services are placed at RSUs along the slotted timeline such that the average response latency for all the requested services can be optimized. The major contributions are summarized as follows.

- A generic approach is put forward for improving the performance of VEC in terms of response latency, with the help of strategic service placement. By formulating the mathematical model, our goal is to minimize the average response latency of requested services along the slotted timeline in this paper.
- Multiple constraints are considered in the optimization problem, especially the overall constraint of energy consumption along the infinite time-slotted horizon. The Lyapunov optimization-based technology is leveraged for converting the long-term energy constraint into per-slot ones. Based on this, an online algorithm is designed to search for the approximate optimal solution in this paper.
- Extensive simulation is conducted for investigating and evaluating the performance of the proposed service placement strategy. Simulation results have revealed that our strategy can achieve better performance compared to other strategies.

The rest of the paper is organized as follows. Some related works are reviewed in Section II. The system model is mathematically formulated in Section III, and the optimization problem is also given in Section IV. An online approach for service placement is put forward in Section V, followed by the simulation evaluation in Section VI. Finally, the conclusion is drawn in Section VII.

## II. RELATED WORK

---

Due to the restricted physical size and energy supply of smart devices (e.g., OBU), there are increasingly strong demands for shifting the task execution from the local to the third party with abundant computational resources. Such a tendency has given birth to several new computing paradigms represented by edge computing. Extensive attention has been paid to this research field, aiming to optimize the performance of task offloading and service outsourcing from different angles [7]. Meanwhile, service placement usually plays an important role in improving the performance of the systems when tasks are offloaded and resources are scheduled in the cloud-edge networks [8–12].

Mobile Edge Computing (MEC), derived from edge computing, brings the computational resources close to the devices and enables tasks to be accomplished with shorter response latency compared to cloud computing. To further enhance the performance of MEC, Li *et al.* [1] jointly optimized task offloading, service placement and resources deployment at the edge of MEC. To improve the quality of service such as response latency, service providers usually deploy services at the network edges. In this context, Panadero *et al.* [13] studied the service placement at the edge server by minimizing the total cost for task offloading which includes the costs on service placement, the throughput of edge server, and the energy consumption. In the meanwhile, the strict requirement for user-perceived delay should be satisfied. They proposed a heuristic algorithm to solve the formulated Mixed Integer Linear Programming problem.

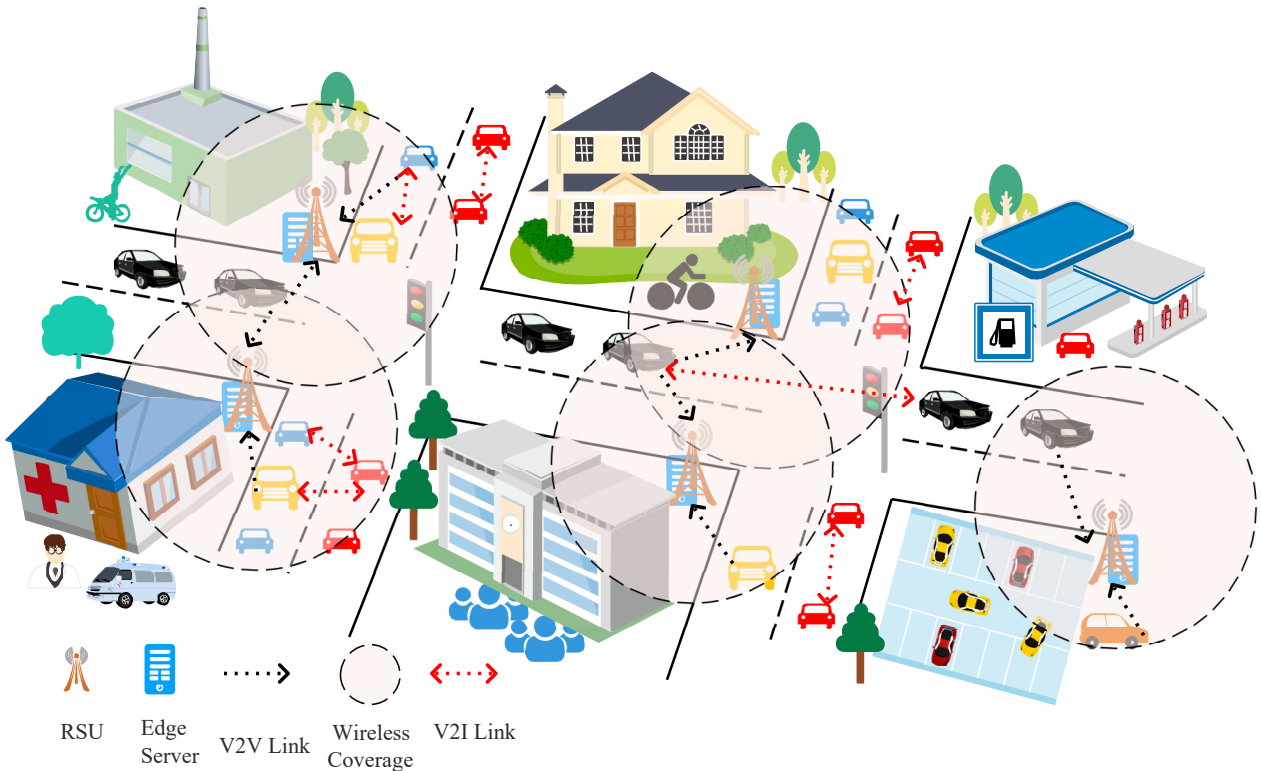
Likewise, Internet of Things (IoT) devices are increasingly demanding services characterized by ultra-low response latencies. In this context, minimization of response latency while meeting multiple constraints has become the main focus in the existing literature. However, current works seldom consider the load-balancing issue which has played a key role in cost-efficient system management. Therefore, an optimization problem with two optimization objectives was proposed in [14], aiming to place IoT services with load balance for improving the quality of service that is defined as a function of deadline violation, service deployment, and unhosted services.

Task execution at the network edges brings benefits to the task requestors, but it also causes serious pressure on the limited computational resources of the edge servers. How to efficiently schedule these

resources for the requestors has become one of the prime challenges. Different from the previous works which try to jointly optimize the service placement and request routing, Yuan *et al.* [15] have considered the non-negligible operating expenses caused by service placement. Specifically, they introduce a two-timescale framework to minimize these costs during frequent cross-cloud service replacement and replication, and further design a greedy algorithm to address the service placement problem. Community networks have attracted lots of attention recently. These networks are usually constructed and managed voluntarily at the edge, and featured by irregular topology and resource heterogeneity. In view of this, Panadero *et al.* [16] proposed a service placement strategy, which can place services in decentralized community networks.

Complicated services can usually be divided into independent components which can be separately accomplished at different edge servers in a cooperative fashion. Previous works usually assumed services are unsplittable or ignored the fact that edge servers can be geographically isolated. Furthermore, the interference caused by services sharing the same physical nodes/links was also ignored, which however can cause long computation delays. Accordingly, Han *et al.* [17] paid attention to the Interference-Aware (IA) service placement while assuming that services can be composed of multiple components in the edge-cloud networks.

Apart from the computational resources, storage and networking resources are both required for data-intensive applications such as machine learning tasks. Furthermore, the service (libraries, databases and codes) placement strategy should be updated over time. Accordingly, a two-time-scale framework is put forward in [18], aiming to jointly optimize service placement and request allocation. Moreover, a polynomial-time service placement algorithm is designed in real-time response to the service requests. Simulation results reveal that their approach achieves 90% of the optimal performance. To realize service migration in dynamic mobile networks, pervasive edge computing is attracting more and more attention recently. Against this background, Ning *et al.* [19] aimed to maximize the utility value of the edge system. Specifically, they consider both storage capacity and response latency in the optimization problem.



**Figure 1.** A VEC application scenario.

Lyapunov optimization is introduced for converting it into instant subproblems. Moreover, two approximation algorithms are used to obtain the average system utility in the future and make decisions on service placement, respectively.

Due to the increasingly unoccupied resources on UAVs, more and more UAVs can act as the edge servers to process the offloaded tasks by the onboard processors [20–22]. Thus, services should be placed beforehand. He *et al.* [20] aimed to jointly optimize the service placement and resource allocation problem in UAV-assisted edge computing, which was modeled as a mixed-integer nonlinear programming problem and divided into two subproblems, which are proven to be submodular and convex, respectively. Thus, a general alternative optimization is adopted to solve this problem.

### III. SYSTEM MODEL

A considered scenario is depicted in Figure 1. Multiple RSUs are scattered in the densely populated area. In addition, their performance is enhanced by deploying edge servers that are rich in computational re-

sources, as denoted in the figure. Thus, services can be provided at the edge of the network. Considering the mobility of vehicles, the strategy for service placement, i.e., which services should be deployed at which RSUs, should be updated periodically. Therefore, we divide the timeline  $\mathcal{T}$  into  $T$  time slots in this paper, and each time slot lasts  $\tau$  seconds, i.e.,  $\mathcal{T} = \{0, \dots, T-1\}$ . The service placement can be updated within each time slot, so as to reflect its dynamics. We denote the set of RSUs by  $\mathcal{N} = \{1, \dots, N\}$ , where  $N$  is the number of RSUs in the considered model. In addition, we further assume that the area considered in this paper can be divided into  $M$  disjoint sub-areas (regions), indexed by  $\mathcal{M}$ . Due to the increasing density of RSU deployment, it is likely that one sub-area or region can be covered by multiple RSUs.

Let  $\mathcal{R}_m(\subset \mathcal{N})$  denote the set of RSUs that can cover the region  $m$ . Thus, vehicles with offloading requests at region  $m$  can resort to RSU  $n \in \mathcal{R}_m$ . Assume that there are  $K$  computing services that can be deployed at RSUs, indexed by  $\mathcal{K} = \{1, \dots, K\}$ . Each service  $k$  can be represented by a 2-tuple of  $(d_k, s_k)$ , where  $d_k$  and  $s_k$  denote the task-input data size and averagely required workload (e.g., the CPU cycles),



respectively. Each RSU  $n$  can be represented by a 2-tuple of  $(D_n, f_n)$ , where  $D_n$  and  $f_n$  denote the caching and processing capabilities of RSU  $n$ , respectively.

**Definition 1. Service Placement Decision.** *Service placement in this paper refers to the deployment of services at the edge server by downloading service-oriented data such as parameters, libraries and databases in advance from the remote data center. Let  $\delta_{k,n}^t$  be a binary variable to indicate whether the service  $k$  is deployed at RSU  $n$  in time slot  $t$ .  $\delta_{k,n}^t = 1$ , if service  $k$  is deployed at RSU  $n$  in time slot  $t$ ; and 0, otherwise. Define  $\delta_n^t = \{\delta_{1,n}^t, \dots, \delta_{K,n}^t\}$  as the service placement decision of RSU  $n$  for all services in  $\mathcal{K}$  in time slot  $t$ , and  $\delta^t = \{\delta_1^t, \dots, \delta_N^t\}$  as the service placement decisions of all RSUs in  $\mathcal{N}$  in time slot  $t$ .*

Note that it is impractical for RSU  $n$  to deploy all the services simultaneously since the storage (caching) capability (i.e.,  $D_n$ ) is limited. In particular, the service placement decisions should meet the following constraint:

$$\sum_{k=1}^K \delta_{k,n}^t d_k \leq D_n, \quad \forall t \in \mathcal{T}, \quad \forall n \in \mathcal{N}, \quad (1)$$

which indicates that the total amount of task-input data of all the services placed at RSU  $n$  should not exceed its caching capability  $D_n$  in each time slot. We investigate the performance of service placement in vehicle edge computing in regions instead of individual vehicles [2], since the statistical information on vehicular offloading requests can be better captured and evaluated in regions than individual vehicles, and this way can better cater to service placement featured by the long-term process. As a result, it is no longer adequate to assume that the number of vehicular offloading requests is constant as most of the existing literature did [3, 5]. In this paper, we assume that the vehicular task arrivals, i.e., the arrivals of the offloading requests generated by vehicles in the corresponding regions, follow a Poisson process. In particular, let  $\lambda_{k,m}^t$  denote the arrival rate of vehicular tasks requesting service  $k$  in the region  $m$  in time slot  $t$ .

Let  $\mathcal{R}_{k,m}^t \subseteq \mathcal{R}_m$  denote the set of RSUs that have deployed service  $k$  in the region  $m$  in time slot  $t$ . Thus, vehicular offloading requests in the region  $m$  can resort to the RSU in  $\mathcal{R}_{k,m}^t$ . Thus, two cases may occur as follows. One is that the set  $\mathcal{R}_{k,m}^t$  is empty. Then,

we assume that the vehicular offloading requests will be forwarded to the cloud center where the tasks are served by corresponding services. The other case is that the set  $\mathcal{R}_{k,m}^t$  is not empty, i.e., there are multiple RSUs that can serve the vehicular offloading requests. For simplicity, the vehicular offloading requests are assumed to be evenly distributed among RSUs in  $\mathcal{R}_{k,m}^t$  [2]. Thus, we define  $w_{k,n}^t$  as the average workloads required by service  $k$  to RSU  $n$  in time slot  $t$ , given as:

$$w_{k,n}^t = \begin{cases} 0 & \bigcup_{m=1}^M \mathcal{R}_{k,m}^t = \emptyset, \\ \sum_{m=1}^M I_{n,m} \frac{s_k \lambda_{k,m}^t}{|\mathcal{R}_{k,m}^t|} & \bigcup_{m=1}^M \mathcal{R}_{k,m}^t \neq \emptyset, \end{cases} \quad (2)$$

where  $\bigcup_{m=1}^M \mathcal{R}_{k,m}^t = \emptyset$  means that there are no regions  $m \in \mathcal{M}$  that have placed service  $k$ .  $I_{n,m}$  is an indicator to denote whether the RSU  $n$  covers the region  $m$ , which can be defined as:

$$I_{n,m} = \begin{cases} 1 & n \in \mathcal{R}_m, \\ 0 & n \notin \mathcal{R}_m. \end{cases} \quad (3)$$

Combining the service placement decision  $\delta_{k,n}^t$ ,  $w_{k,n}^t$  can be further rewritten as:

$$w_{k,n}^t = \delta_{k,n}^t \sum_{m=1}^M I_{n,m} \frac{s_k \lambda_{k,m}^t}{|\mathcal{R}_{k,m}^t|}, \quad \forall k \in \mathcal{K}, \forall n \in \mathcal{N}. \quad (4)$$

Service placement at the edge server in advance aims to reduce the response delay and further improve the Quality of Service (QoS) and Quality of Experience (QoE). However, it shall be noted that the service deployment and placement at the edge, e.g., by downloading task-input data, libraries and databases pertaining to the corresponding services from the data center, will also incur costs such as energy consumption at the edge.

### 3.1 Response Delay Model

The average response delay for the tasks processed at the edge mainly consists of the queueing delay and computation delay, when the corresponding services have already been placed at the edge server in advance. We have considered the queueing delay in the paper,

owing to the fact that the computational resources at RSUs cannot process all the tasks in parallel, especially when the number of vehicular offloading requests is huge. Furthermore, the queuing delay is one of the causes that incur long response latency.

Since the vehicular task arrival of each service type follows a Poisson process with rate  $\lambda_{k,m}^t$  and the offloading requests are evenly distributed among RSUs in  $\mathcal{R}_{k,m}^t$ , the overall task arrival at RSU  $n$  for service  $k$  is also a Poisson process, with the rate  $\lambda_{k,n}^t$  given as:

$$\lambda_{k,n}^t = \sum_{m=1}^M I_{n,m} \frac{\lambda_{k,m}^t}{|\mathcal{R}_{k,m}^t|}, \quad \forall k \in \mathcal{K}, \forall n \in \mathcal{N}. \quad (5)$$

The service time of one task of service type  $k$  at RSU  $n$  can be expressed as  $s_k/f_n$ . Accordingly, the service rate  $\mu_{k,n}$  is:

$$\mu_{k,n} = \frac{f_n}{s_k}, \quad \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \quad (6)$$

where  $f_n$  denotes the processing frequency of RSU  $n$ . Specifically, we evaluate and analyze the response delay by modeling this process as an M/M/1 queue. Thus the response delay, i.e., the average sojourn time (waiting time and service time) for service  $k$  at RSU  $n$  in time slot  $t$ , is given as:

$$rd_{k,n}^t = \frac{s_k}{f_n - s_k \sum_{m=1}^M \frac{I_{n,m} \lambda_{k,m}^t}{|\mathcal{R}_{k,m}^t|}}. \quad (7)$$

The above equation implies that service  $k$  has been placed at RSU  $n$ . On the other hand, it is also possible that service  $k$  is not placed at RSU  $n$ . In this case, the offloading request will be served at the cloud center. Due to the sufficient computing resources and powerful computing capabilities at the cloud center, we thus assume that the calculation delay at the cloud center can be ignored. The response delay, denoted by  $rd_{k,c}^t$ , is mainly caused by the transmission delay when data and execution results are transmitted over the backbone network using the backhaul links. Such transmission delay can be estimated and predicted based on historical statistical data and highly efficient machine learning approaches such as auto-regression analysis. Hence, the average response delay for service  $k$  in time

slot  $t$   $rd_k^t$  is:

$$rd_k^t = \frac{1}{N} \sum_{n=1}^N \delta_{k,n}^t rd_{k,n}^t + (1 - \delta_{k,n}^t) rd_{k,c}^t. \quad (8)$$

### 3.2 Energy Consumption Model

As mentioned earlier, the service placement at the edge in advance incurs negligible energy consumption, including static energy consumption and dynamic energy consumption. The static energy, denoted by  $\alpha_n$ , is consumed for maintaining virtual resources allocated to service  $k$ , which is independent of the workload of the service. The dynamic energy  $\beta_{k,n}^t$  is consumed for performing the service at run time. Thus, the total energy consumption for service  $k$  at RSU  $n$  in time slot  $t$  (i.e., the static energy consumption plus the dynamic energy consumption) can be expressed as:

$$\begin{aligned} e_{k,n}^t &= \alpha_n + \varsigma \eta w_{k,n}^t (f_n)^2 \\ &= \alpha_n + \varsigma \eta \delta_{k,n}^t (f_n)^2 \sum_{m=1}^M I_{n,m} \frac{s_k \lambda_{k,m}^t}{|\mathcal{R}_{k,m}^t|}, \end{aligned} \quad (9)$$

where the second part on the right-hand side represents the dynamic energy consumption  $\beta_{k,n}^t$ ,  $\varsigma$  represents the effectively switched capacitance coefficient, and  $\eta$  is the number of cycles needed to perform one task-input bit from  $n$ . The energy consumption for computation undertaking at the edge server is supposed to be lower than the given threshold, so the residual energy can be kept for other purposes such as vehicle-to-infrastructure and pedestrian-to-infrastructure communications. Thus, the following constraint should be satisfied:

$$\sum_{k=1}^K e_{k,n}^t \leq e_{max,n}^t, \quad \forall n \in \mathcal{N}, \forall t \in \mathcal{T}, \quad (10)$$

where  $e_{max,n}^t$  is the maximal energy consumption that RSU  $n$  can tolerate to perform the tasks. Therefore, for arbitrary time slot  $t$ , the total energy consumption for undertaking all the services at RSU  $n$  should not exceed this energy threshold. For easy discussion, we assume that the energy threshold is the same for each time slot at RSU  $n$ .

## IV. PROBLEM FORMULATION

Our goal in this paper is to minimize the overall response delay of services along the infinite time slots, by adjusting the strategy of service placement at RSUs within each time slot. In particular, our optimization problem can be formulated as:

$$(\mathcal{P}1) \min_{\delta^t, \forall t} \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{k=1}^K r d_k^t,$$

$$\text{s.t. : } \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{k=1}^K e_{k,n}^t \leq \mathcal{E}_n, \quad \forall n \in \mathcal{N}, \quad (11)$$

$$f_{n,\min} \leq f_n \leq f_{n,\max}, \quad \forall n \in \mathcal{N}, \quad (12)$$

$$I_{n,m} \in \{0, 1\} \quad \forall n \in \mathcal{N}, \forall m \in \mathcal{M}, \quad (13)$$

$$\delta_{k,n}^t \in \{0, 1\}, \quad \forall k \in \mathcal{K}, \forall n \in \mathcal{N}, \forall t \in \mathcal{T}, \quad (14)$$

$$\text{Constraints (1), (10),} \quad (15)$$

where (11) represents that the long-term energy consumption at RSU  $n$  should not exceed the threshold  $\mathcal{E}_n$ , for the reason that the evaluation of service placement strategy needs a long process as mentioned earlier. The computational resources allocated for service deployment and execution should be bounded as shown in (12), owing to the limited computing capabilities at the edge server in contrast to the cloud center. The number of services that are placed is restricted by the limited storage capacity of the edge server. In addition, per-slot energy consumption for RSU  $n$  is also bounded by  $e_{\max,n}^t$  as shown in (10).

The difficulties in addressing the problem  $\mathcal{P}1$  are summarized as follows. First, the entire information, such as the offloading requests in all time slots, is required for optimally solving this problem. However, it is difficult to obtain this information about the future time slots in the current time slot. In addition,  $\mathcal{P}1$  is essentially an integer linear programming that is computationally prohibitive even if the entire information is acquired a priori. Last but not least, it is quite difficult to handle the constraint (11), since it similarly requires the entire information in all the time slots. Accordingly, an online approach is needed for addressing this problem, such that the response delay can be optimized by placing services at the edge within each time slot.

## V. ONLINE ALGORITHM FOR SERVICE PLACEMENT

To address the above challenges, we in this section present an online algorithm, named Online service Placement for vehicular Edge Computing (OPEC). Specifically, the Lyapunov optimization technology is introduced for converting the infinite time-slotted (global) energy consumption constraint (11) into respective constraints within each time slot. Thus, the optimization for  $\mathcal{P}1$  can be solved approximately within each time slot.

### 5.1 Lyapunov-Based Online Algorithm

To covert the global constraint into local ones,  $n$  energy migration queues  $Q_n(n \in \mathcal{N})$  are constructed, which can be recursively defined as:

$$Q_n(t+1) = \max[Q_n(t) - \mathcal{E}_n, 0] + \sum_{k=1}^K e_{k,n}^t, \quad \forall n \in \mathcal{N}, \quad (16)$$

where  $Q_n(t)$  is a numerical value to denote the queue backlog of RSU  $n$  in time slot  $t$ . Actually,  $Q_n(t)$  can indicate how far the energy consumption in the current time slot deviates from the global energy constraint  $\mathcal{E}_n$ , e.g., a larger value of  $Q_n(t)$  for RSU  $n$  always denotes a sharper deviation from  $\mathcal{E}_n$ . Note that we assume that  $Q_n(0) = 0$  holds,  $\forall n \in \mathcal{N}$ .

**Lemma 1.** For an arbitrary RSU  $n \in \mathcal{N}$ , given the corresponding queue backlog  $Q_n(t)$  in time slot  $t$ , the following inequality always holds

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n] \leq \frac{1}{T} \mathbb{E}[Q_n(T)].$$

*Proof.* The lemma can be proven based on the following two cases.

- Case 1:  $Q_n(t) \geq \mathcal{E}_n$ . In this case, we have  $Q_n(t+1) = \max[Q_n(t) - \mathcal{E}_n, 0] + \sum_{k=1}^K e_{k,n}^t = Q_n(t) - \mathcal{E}_n + \sum_{k=1}^K e_{k,n}^t$ , so we can obtain  $\sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n = Q_n(t+1) - Q_n(t)$ .
- Case 2:  $Q_n(t) < \mathcal{E}_n$ . In this case,  $Q_n(t+1) = \max[Q_n(t) - \mathcal{E}_n, 0] + \sum_{k=1}^K e_{k,n}^t = \sum_{k=1}^K e_{k,n}^t$  and  $\sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n = Q_n(t+1) - Q_n(t)$ . As a result,

$\sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n \leq Q_n(t+1) - Q_n(t)$  holds,  $\forall t \in \mathcal{T}$ . Take the expectation of this inequality and then summarize it over  $\forall t \in \mathcal{T}$ , we have  $\sum_{t=0}^{T-1} \mathbb{E} \left[ \sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n \right] \leq \sum_{t=0}^{T-1} \mathbb{E} [Q_n(t+1) - Q_n(t)] = \mathbb{E} [Q_n(T)]$ . Thus,  $\forall n \in \mathcal{N}$ , we can obtain  $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n \right] \leq \frac{1}{T} \mathbb{E} [Q_n(T)]$ .

Define the vector  $\mathbf{Q}(t) = \{Q_1(t), \dots, Q_N(t)\}$  as the queue backlogs of  $N$  RSUs in time slot  $t$ . Then, the Lyapunov function with regards to (w.r.t.)  $\mathbf{Q}(t)$  is defined as:  $\mathcal{F}(\mathbf{Q}(t)) = \frac{1}{2} \sum_{n=1}^N Q_n^2(t)$ .

The Lyapunov function is actually the function of the queue backlogs of all the RSU  $\mathcal{N}$  in time slot  $t$ . We can further define *Lyapunov drift* as the increment of the queue backlogs in two consecutive time slots, given as:

$$\Delta(\mathbf{Q}(t)) \triangleq \mathbb{E}[\mathcal{F}(\mathbf{Q}(t+1)) - \mathcal{F}(\mathbf{Q}(t)) | \mathbf{Q}(t)]. \quad (17)$$

The Lyapunov drift can actually indicate the fluctuation of two consecutive states. From the perspective of system stability, the smaller the Lyapunov drift, the better the system stability. However, minimizing the Lyapunov drift faces a similar challenge, since it also needs the information in future time slots. For example, the optimization of  $\Delta(\mathbf{Q}(t))$  needs the future information  $\mathcal{F}(\mathbf{Q}(t+1))$ . To avoid involving future information, we can actually relax it by determining its upper bound. In the next, we will show how to find the upper bound of the Lyapunov drift according to the following lemma [23].

**Lemma 2.** *Given four non-negative real numbers  $A$ ,  $B$ ,  $C$  and  $m$ , they satisfy  $C = \max\{B - m, 0\} + A$ , then  $C^2 \leq A^2 + B^2 + m^2 - 2B(m - A)$ .*

Based on this Lemma, we have

$$\begin{aligned} \Delta(\mathbf{Q}(t)) &= \mathbb{E}[\mathcal{F}(\mathbf{Q}(t+1)) - \mathcal{F}(\mathbf{Q}(t)) | \mathbf{Q}(t)] \\ &= \frac{1}{2} \mathbb{E} \left[ \sum_{n=1}^N [\max\{Q_n(t) - \mathcal{E}_n, 0\} + \sum_{k=1}^K e_{k,n}^t]^2 \right. \\ &\quad \left. - \sum_{n=1}^N Q_n^2(t) | \mathbf{Q}(t) \right] \\ &\leq \frac{1}{2} \mathbb{E} \left[ \sum_{n=1}^N [\mathcal{E}_n^2 + (\sum_{k=1}^K e_{k,n}^t)^2] \right] \end{aligned}$$

$$\begin{aligned} &+ 2Q_n(t) (\sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n) | \mathbf{Q}(t) \\ &= D - \sum_{n=1}^N \mathcal{E}_n Q_n(t) + \mathbb{E} \left[ \sum_{n=1}^N Q_n(t) \sum_{k=1}^K e_{k,n}^t | \mathbf{Q}(t) \right], \end{aligned}$$

where

$$\begin{aligned} D &= \frac{1}{2} \sum_{n=1}^N \mathcal{E}_n^2 + \frac{1}{2} \mathbb{E} \left[ \sum_{n=1}^N \left( \sum_{k=1}^K e_{k,n}^t \right)^2 | \mathbf{Q}(t) \right] \\ &\leq \frac{1}{2} \sum_{n=1}^N \mathcal{E}_n^2 + \frac{1}{2} \mathbb{E} \left[ \sum_{n=1}^N (e_{max,n}^t)^2 | \mathbf{Q}(t) \right] \\ &= \frac{1}{2} \sum_{n=1}^N [\mathcal{E}_n^2 + (e_{max,n}^t)^2] \triangleq B. \end{aligned}$$

Therefore, the upper bound of the Lyapunov drift can be given as:

$$\begin{aligned} \Delta(\mathbf{Q}(t)) &\leq B - \sum_{n=1}^N \mathcal{E}_n Q_n(t) \\ &\quad + \mathbb{E} \left[ \sum_{n=1}^N Q_n(t) \sum_{k=1}^K e_{k,n}^t | \mathbf{Q}(t) \right]. \quad (18) \end{aligned}$$

It is noticeable that the upper bound of  $\Delta(\mathbf{Q}(t))$  no longer depends upon the information in future time slots such as  $\mathbf{Q}(t+1)$ . Then, we introduce the *drift-plus-penalty* term that integrates the Lyapunov drift into the optimization objective per time slot, given as:

$$\Theta(t) = \Delta(\mathbf{Q}(t)) + V \mathbb{E} \left[ \sum_{k=1}^K r d_k^t | \mathbf{Q}(t) \right], \quad (19)$$

where  $V$  is a positive numeric value to denote the preference between energy consumption and response delay. Combining (18), we have:

$$\begin{aligned} \Theta(t) &\leq B - \sum_{n=1}^N \mathcal{E}_n Q_n(t) \\ &\quad + \mathbb{E} \left[ \sum_{n=1}^N Q_n(t) \sum_{k=1}^K e_{k,n}^t + V \sum_{k=1}^K r d_k^t | \mathbf{Q}(t) \right]. \quad (20) \end{aligned}$$

We now turn our attention to the optimization of the supremum of the drift-plus-penalty term, i.e., the



right-hand side of (20). Since  $B$ ,  $\mathcal{E}_n$  and  $Q_n(t)$  are independent of the service placement decisions and thus can be obtained in time slot  $t$ , to optimize the supremum of the drift-plus-penalty is equivalent to the optimization of the following problem:

$$(\mathcal{P}2) \quad \min_{\forall t, \delta^t} \left\{ \sum_{n=1}^N Q_n(t) \sum_{k=1}^K e_{k,n}^t + V \sum_{k=1}^K rd_k^t \right\},$$

s.t. : (12), (13), (14), (15). (21)

Therefore, we can solve the problem  $\mathcal{P}1$  approximately by solving the problem  $\mathcal{P}2$ . Compared to  $\mathcal{P}1$ , it is relatively easier to solve  $\mathcal{P}2$ , because the time slot spanned energy constraint has been converted into per-slot ones. An important issue is naturally posed, i.e., what is the degree of approximation of the solution of  $\mathcal{P}2$  towards that of  $\mathcal{P}1$ ?

**Theorem 1.** *If there exists a solution to the problem  $\mathcal{P}1$ , the service placement profile obtained by solving  $\mathcal{P}2$  over time slots  $t \in \{0, \dots, T-1\}$  can also approximately optimize the problem  $\mathcal{P}1$  and the optimality gap which denotes the difference between the two solutions is bounded by  $O(1/V)$ .*

*Proof.* Since there is a solution to  $\mathcal{P}1$ , there exists a service placement profile  $\delta^{t*}$  in time slot  $t$ , and the corresponding global optimum  $l^*$ , which respectively satisfy 1)  $\sum_{k=1}^K e_{k,n}^t \leq \mathcal{E}_n, \forall n \in \mathcal{N}$  and 2)  $l^* = \min_{\delta^t, \forall t} \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \sum_{k=1}^K rd_k^t$ , based on [24].

Accordingly, we have:

$$\begin{aligned} \Theta(t) &= \Delta(\mathbf{Q}(t)) + V \mathbb{E} \left[ \sum_{k=1}^K rd_k^t | \mathbf{Q}(t) \right] \\ &\leq B - \sum_{n=1}^N \mathcal{E}_n Q_n(t) \\ &\quad + \mathbb{E} \left[ \sum_{n=1}^N Q_n(t) \sum_{k=1}^K e_{k,n}^t + V \sum_{k=1}^K rd_k^t | \mathbf{Q}(t) \right] \\ &= B + \mathbb{E} \left[ \sum_{n=1}^N Q_n(t) \left( \sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n \right) | \mathbf{Q}(t) \right] \\ &\quad + V \mathbb{E} \left[ \sum_{k=1}^K rd_k^t | \mathbf{Q}(t) \right] \\ &\stackrel{\ddagger}{\leq} B + Vl^*. \end{aligned} \quad (22)$$

Since  $\Theta(t) \leq B + \mathbb{E}[\sum_{n=1}^N Q_n(t)(\sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n) | \mathbf{Q}(t)] + V \mathbb{E}[\sum_{k=1}^K rd_k^t | \mathbf{Q}(t)]$ , for  $\forall t \in \{0, \dots, T-1\}$  and valid service placement profile  $\delta^t$ , the inequality still holds when substituting  $\delta^t$  by  $\delta^{t*}$ . Owing to  $\sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n \leq 0$  for  $\delta^{t*}$ , the inequality ( $\ddagger$ ) holds.

Take the expectation of (22) and add it over the time slots  $t \in \{0, \dots, T-1\}$ , i.e.,

$$\begin{aligned} &\sum_{t=0}^{T-1} \mathbb{E}[\Delta(\mathbf{Q}(t)) + V \mathbb{E}[\sum_{k=1}^K rd_k^t | \mathbf{Q}(t)]] \\ &= \sum_{t=0}^{T-1} \mathbb{E}[\mathbb{E}[\mathcal{F}(\mathbf{Q}(t+1)) - \mathcal{F}(\mathbf{Q}(t)) | \mathbf{Q}(t)] \\ &\quad + V \mathbb{E}[\sum_{k=1}^K rd_k^t | \mathbf{Q}(t)]] \\ &= \sum_{t=0}^{T-1} \mathbb{E}[\mathcal{F}(\mathbf{Q}(t+1)) - \mathcal{F}(\mathbf{Q}(t))] + V \mathbb{E}[\sum_{k=1}^K rd_k^t] \\ &= \mathbb{E}[\mathcal{F}(\mathbf{Q}(T)) - \mathcal{F}(\mathbf{Q}(0))] + \sum_{t=0}^{T-1} V \mathbb{E}[\sum_{k=1}^K rd_k^t] \\ &= \mathbb{E}[\mathcal{F}(\mathbf{Q}(T))] + \sum_{t=0}^{T-1} V \mathbb{E}[\sum_{k=1}^K rd_k^t] \\ &\leq \sum_{t=0}^{T-1} \mathbb{E}[B + Vl^*] = (B + Vl^*) * T. \end{aligned}$$

Since  $\mathbb{E}[\mathcal{F}(\mathbf{Q}(T))] \geq 0$ , we can have  $\sum_{t=0}^{T-1} V \mathbb{E}[\sum_{k=1}^K rd_k^t] \leq (B + Vl^*) * T$ . Thus, we obtain  $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\sum_{k=1}^K rd_k^t] \leq l^* + \frac{B}{V}$ . Since  $\sum_{k=1}^K rd_k^t$  is the expected response time for all the services in time slot  $t$ , and according to the law of iterated expectations,  $\mathbb{E}[\sum_{k=1}^K rd_k^t] = \sum_{k=1}^K rd_k^t$ . Therefore, we have  $\frac{1}{T} \sum_{t=0}^{T-1} \sum_{k=1}^K rd_k^t \leq l^* + \frac{B}{V}$ .

Thus, the service placement profile obtained by solving  $\mathcal{P}2$  can infinitely approximate the optimization of  $\mathcal{P}1$  by adjusting the variable  $V$ . Let  $\mathcal{G}(V)$  denote the optimality gap which denotes the difference between the two solutions. The optimal value by solving the supremum of drift-plus-penalty term is  $l^* + B/V$ , so  $\mathcal{G}(V) = B/V$  and  $O(\mathcal{G}(V)) = O(1/V)$ . Hence, the optimality gap is bounded by  $O(1/V)$ .

The above analysis can guarantee that the service placement profile obtained by solving  $\mathcal{P}2$  can approximately solve  $\mathcal{P}1$ . However, another issue that needs

to be addressed is whether the service placement decisions obtained from  $\mathcal{P}2$  do not violate the long-term constraint (11) for  $\mathcal{P}1$ .

**Theorem 2.** *The service placement decisions obtained from solving  $\mathcal{P}2$  over time slot  $t \in \{0, \dots, T-1\}$  always satisfy (11).*

*Proof.* Let  $\delta^{t*}$  denote the service placement decisions of all RSUs  $\mathcal{N}$  in time slot  $t$  obtained by solving problem  $\mathcal{P}2$ , i.e., the supremum of the drift-plus-penalty term. In other words, when services are placed based on  $\delta^{t*}$ , the supremum of the drift-plus-penalty term can achieve the minimum. For the sake of clear expression and easy discussion, let  $Y(\delta^t) = \sum_{k=1}^K rd_k^t$  and  $Z_n(\delta^t) = \sum_{k=1}^K e_{k,n}^t$ . We have:

$$\begin{aligned} \Theta(t) &= \Delta(\mathbf{Q}(t)) + V\mathbb{E}[Y(\delta^{t*})|\mathbf{Q}(t)] \\ &\leq B - \sum_{n=1}^N \mathcal{E}_n Q_n(t) \\ &\quad + \mathbb{E}\left[\sum_{n=1}^N Q_n(t) Z_n(\delta^{t\dagger}) + VY(\delta^{t\dagger})|\mathbf{Q}(t)\right] \\ &\stackrel{\S}{\leq} B - \sum_{n=1}^N \mathcal{E}_n Q_n(t) + \sum_{n=1}^N Q_n(t)(\mathcal{E}_n - \epsilon_n) + VL_t^\dagger \\ &= B - \sum_{n=1}^N \epsilon_n Q_n(t) + VL_t^\dagger, \end{aligned}$$

where  $\delta^{t\dagger}$  is an arbitrary valid service placement profile for all RSUs  $\mathcal{N}$  in time slot  $t$  except  $\delta^{t*}$  and  $L_t^\dagger$  is the corresponding response delay. Due to  $\mathbb{E}[Z_n(\delta^{t*})] = Z_n(\delta^{t*}) \leq \mathcal{E}_n$ , there exists a small enough  $\epsilon_n > 0$  that satisfies  $Z_n(\delta^{t*}) \leq \mathcal{E}_n - \epsilon_n$ ,  $\forall n \in \mathcal{N}$ , so the inequality (§) holds.

Take the expectations of the above inequality,  $\mathbb{E}\{\Theta(t)\} \leq B - \mathbb{E}[\sum_{n=1}^N \epsilon_n Q_n(t)] + VL_t^\dagger$ , i.e.,

$$\begin{aligned} &\mathbb{E}[\mathcal{F}(\mathbf{Q}(t+1))] - \mathbb{E}[\mathcal{F}(\mathbf{Q}(t))] + V\mathbb{E}[Y(\delta^{t*})] \\ &\leq B + VL_t^\dagger - \mathbb{E}\left[\sum_{n=1}^N \epsilon_n Q_n(t)\right]. \end{aligned}$$

By summing this inequality over  $t \in \{0, \dots, T-1\}$ , we have

$$\mathbb{E}[\mathcal{F}(\mathbf{Q}(T))] - \mathbb{E}[\mathcal{F}(\mathbf{Q}(0))] + V \sum_{t=0}^{T-1} \mathbb{E}[Y(\delta^{t*})]$$

$$\begin{aligned} &\leq BT + VL_t^\dagger T - \mathbb{E}\left[\sum_{t=0}^{T-1} \sum_{n=1}^N \epsilon_n Q_n(t)\right] \\ &\leq BT + VL_t^\dagger T = B'T, \end{aligned}$$

where  $B' = B + VL_t^\dagger$ . Since  $\mathcal{F}(\mathbf{Q}(T)) = \frac{1}{2} \sum_{n=1}^N Q_n^2(T)$ , we have  $\mathbb{E}[\sum_{n=1}^N Q_n^2(T)] \leq 2B'T$ . Then, we have  $\mathbb{E}^2\left[\sum_{n=1}^N Q_n(T)\right] \leq K\mathbb{E}\left[\sum_{n=1}^N Q_n^2(T)\right] \leq 2B'T$ . Therefore, we have  $\mathbb{E}\left[\sum_{n=1}^N Q_n(T)\right]/T \leq \sqrt{2B'/T}$ . Due to  $Q_n(T) \geq 0$  for  $\forall n \in \mathcal{N}$ ,  $\mathbb{E}[Q_n(T)]/T \leq \sqrt{2B'/T}$  holds, for  $\forall n \in \mathcal{N}$ . Given Lemma 1, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[\sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n\right] \leq \frac{1}{T} \mathbb{E}[Q_n(T)] \leq \sqrt{\frac{2B'}{T}}.$$

After taking the limit on both sides, we have

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\left[\sum_{k=1}^K e_{k,n}^t - \mathcal{E}_n\right] \leq \lim_{T \rightarrow \infty} \sqrt{\frac{2B'}{T}} = 0,$$

i.e.,  $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{k=1}^K e_{k,n}^t \leq \mathcal{E}_n$ ,  $\forall n \in \mathcal{N}$ .

Therefore, the following conclusions hold 1) the queue backlog is mean rate stable; 2) the service placement profile for  $\mathcal{P}2$  can solve  $\mathcal{P}1$  without violating the constraint (11).

## 5.2 OPEC Algorithm Description

**Algorithm 1.** *Online service placement for vehicular edge computing (OPEC).*

---

**Input:**  $Q_n(0), \mathcal{E}_n, D_n, N, M, T$   
**Output:** Service placement decisions

- 1:  $\mathcal{F}(\mathbf{Q}(0)) = 0$
- 2: Calculate  $\mu_{k,n}$  based on (6),  $\forall k \in \mathcal{K}, \forall n \in \mathcal{N}$
- 3: **for**  $t = 0$  to  $T - 1$  **do**
- 4:   Observe and record  $\lambda_{k,m}^t, \mathcal{R}_{k,m}^t, \forall k \in \mathcal{K}, \forall m \in \mathcal{M}$
- 5:   Calculate  $\lambda_{k,n}^t$  based on (5),  $\forall k \in \mathcal{K}, \forall n \in \mathcal{N}$
- 6:   Set  $e_{max,n}^t, \forall n \in \mathcal{N}$
- 7:   Predict and obtain  $rd_{k,c}^t, \forall k \in \mathcal{K}$
- 8:   Obtain  $\delta^{t*}$  by solving problem  $\mathcal{P}2$
- 9:   Calculate  $e_{k,n}^t$  based on (9),  $\forall k \in \mathcal{K}, \forall n \in \mathcal{N}$
- 10:   Update  $Q_n(t+1)$  based on (16),  $\forall n \in \mathcal{N}$
- 11: **end for**

---

The algorithm OPEC, which leverages the Lya-

punov technology to convert the time slot spanned energy consumption, is shown in Algorithm 1. This algorithm aims to obtain the optimal service placement decisions for optimizing problem  $\mathcal{P}1$ . To this end, the service placement decisions can be obtained by solving the problem  $\mathcal{P}2$  along each time slot. In the beginning, some variables should be initialized such as  $\mathcal{F}(Q(0))$  and  $\mu_{k,n}$ . Within each time slot, some variables are observed and calculated such as  $\lambda_{k,m}^t$ ,  $\mathcal{R}_{k,m}^t$ , and  $\lambda_{k,n}^t$ . Then the approximately optimum service placement decision is obtained by solving problem  $\mathcal{P}2$  (line 8). After that, the algorithm calculates the energy consumption based on the service placement decisions (line 9) and further updates the queue backlog of each queue (line 10).

### 5.3 Algorithm Design for $\mathcal{P}2$

In algorithm OPEC, there is still a pending issue that needs to be addressed, i.e., the optimization of  $\mathcal{P}2$ , so in this subsection, we will discuss how to obtain the decision placement decision in each time slot by solving the problem  $\mathcal{P}2$  efficiently, i.e.,  $\delta^{t*} = \arg \min_{\delta^t, \forall t} \{ \sum_{n=1}^N Q_n(t) \sum_{k=1}^K e_{k,n}^t + V \sum_{k=1}^K r d_k^t \}$ .

**Lemma 3.** *The optimization problem  $\mathcal{P}2$  is NP-Hard.*

*Proof.* It is well-known that the multiple knapsack problem (MKP) is an NP-hard problem [25], which can be formally described as follows. Given a set of knapsacks  $\mathcal{G}$  and items  $\mathcal{I}$ , each knapsack  $k \in \mathcal{G}$  can pack finite items based on its knapsack capacity  $a_k$  and each item  $i \in \mathcal{I}$  has a size  $s_i$  and profit  $p_i$ . MKP aims to find a way to assign items into the knapsacks such that the total profits can be maximized, i.e.,

$$(\mathcal{P}3) \quad \max \sum_{k \in \mathcal{K}} \sum_{i \in \mathcal{I}} x_{k,i} p_i,$$

$$\text{s.t. :} \quad \sum_{i \in \mathcal{I}} x_{k,i} s_i \leq a_k \quad \forall k \in \mathcal{G}, \quad (23)$$

$$x_{k,i} \in \{0, 1\} \quad \forall k \in \mathcal{G}, \forall i \in \mathcal{I}. \quad (24)$$

Next, we will prove that our problem is actually an instance of MKP, i.e.,  $\mathcal{P}2$  is equivalent to  $\mathcal{P}3$ . First, it is easy to transform  $\mathcal{P}2$  into a maximization problem, e.g., by the inverse operation. Then, given our

system model, the set of RSUs  $\mathcal{N}$  and the set of services  $\mathcal{K}$  correspond to the set of knapsacks  $\mathcal{G}$  and items  $\mathcal{I}$ , respectively. Each service has the required amount of task-input data (e.g.,  $d_k$ ) which corresponds to the item size  $s_i$ . Each RSU has a storage capacity  $D_n$  that corresponds to the backpack capacity  $a_k$ . The total amount of task-input data for all the tasks assigned to the RSU should not exceed its storage capacity. Such a constraint is actually the knapsack constraint in (23). The objective function can be considered as the profit of each service in  $\mathcal{P}2$ . Therefore,  $\mathcal{P}2$  is equivalent to  $\mathcal{P}3$  where  $K$  items are packed to  $N$  RSUs to optimize the response delay. Accordingly, the optimization problem  $\mathcal{P}2$  is NP-Hard.

The searching space in the potential solution domain is up to  $K^N$ , which thus prohibits the exhaustive search in reality, especially considering the strict latency requirements of vehicular services. In addition, iteration-based evolutionary algorithms have similar drawbacks in handling the strict latency requirements. Indeed, the time taken to make decisions on service placement at RSUs should be as fast as possible. Thus, we in this paper put forward a greedy algorithm for this problem to cater to such requirements in the next.

Actually, several factors could affect the performance of the service placement decision in pursuit of energy consumption and response latency optimization. For instance, the request rate of service usually indicates its popularity. Intuitively, placing the most popular services at RSUs can reduce the response delay to the most extent. Further to this, services with larger sizes of task-input data tend to be placed at RSUs, since they actually consume more network resources than those with smaller task-input data sizes. Next, services that require less computational resources tend to be placed, since services with larger computational demands usually incur relatively long computation delays and large energy consumption. Last but not least, the value of  $V$  can also affect the objective value of  $\mathcal{P}2$ .  $V$  is used to adjust the tradeoff between response latency and energy consumption. Usually, a larger value of  $V$  indicates that the VEC system values the response latency more than the energy consumption, while a smaller value of  $V$  indicates that the system values energy consumption more due to the limited energy supply. Specifically, we use the average service requests arriving at RSU  $n$  for service  $k$  denoted by  $p_{k,n}^t$  as the popularity of

service  $k$ , and  $p_{k,n}^t$  is defined as:

$$p_{k,n}^t = \sum_{m=1}^M I_{n,m} \frac{\lambda_{k,m}^t}{|\mathcal{R}_m|}, \quad \forall k \in \mathcal{K}, \forall n \in \mathcal{N}. \quad (25)$$

Based on these observations, we design a utility function that incorporates the above three factors into one single value to evaluate each service within each time slot. In particular, the utility function of service  $k$  at RSU  $n$  in time slot  $t$  can be defined as:

$$\begin{aligned} \mathcal{U}_{k,n}(t) = & w_v \frac{V_{max} - V}{V_{max} - V_{min}} + w_p \frac{p_{k,n}^t - \min_{l \in \mathcal{K}}(p_{l,n}^t)}{\max_{l \in \mathcal{K}}(p_{l,n}^t) - \min_{l \in \mathcal{K}}(p_{l,n}^t)} \\ & + w_d \frac{d_k - \min_{l \in \mathcal{K}}(d_l)}{\max_{l \in \mathcal{K}}(d_l) - \min_{l \in \mathcal{K}}(d_l)} + w_s \frac{\max_{l \in \mathcal{K}}(s_l) - s_k}{\max_{l \in \mathcal{K}}(s_l) - \min_{l \in \mathcal{K}}(s_l)}, \end{aligned} \quad (26)$$

where  $w_p, w_d, w_s$  and  $w_v$  are numerical values to denote the preferences towards the four factors which satisfy  $w_p + w_d + w_s + w_v = 1$ . The two functions  $\max(\cdot)$  and  $\min(\cdot)$  denote the maximal and minimal values, respectively.  $V_{max}$  and  $V_{min}$  are the upper and lower bounds of  $V$ , respectively. Therefore, the service with a larger utility value is more worth being placed at RSUs.

A greedy algorithm GAPO for solving  $\mathcal{P}2$  is put forward as depicted in Algorithm 2. As illustrated in our system model, RSUs are inter-connected with each other, so the information collected from regions such as  $\lambda_{k,m}^t$  can be timely disseminated and shared among them. Some initializations should be accomplished such as  $\delta^t = \mathbf{0}$ . Additionally, the algorithm initializes one Max-Heap  $h_n$  for each RSU  $n$  within each time slot. In the next, each RSU can make its own service placement decisions independently. Specifically, RSU  $n$  calculates the utility value of each service based on the gathered information. The heuristic rule is that services with larger utility values tend to be placed at RSUs. To achieve this goal, the algorithm pushes the pair of service index and the corresponding utility value (i.e.,  $(k, \mathcal{U}_{k,n}^t)$ ) into the heap  $h_n$ . Thus, a service with a larger utility value is always retrieved first when the pop operation over  $h_n$  is implemented. Then  $n$  checks whether there is enough storage capacity for the current service popped from  $h_n$ . If the storage ca-

capacity is sufficient, the current service is placed and the service placement decision is updated (line 15). Repeat the procedure until  $h_n$  is empty or RSU  $n$  has insufficient storage capacity.

---

**Algorithm 2.** Greedy algorithm for per slot optimization (GAPO).

---

**Input:**  $\mathcal{K}, \mathcal{N}, V, \mathcal{E}_n, D_n, w_p, w_d, w_s, t$

**Output:**  $\delta^t$

- 1: Record and share  $\lambda_{k,m}^t$  among RSUs  $\mathcal{N}$ ,  $\forall k \in \mathcal{K}$ ,  $\forall m \in \mathcal{M}$
  - 2: Initialize  $\delta^t$  with zeros
  - 3: Initialize a heap  $h_n$  for each  $n \in \mathcal{N}$
  - 4: **for** each  $n$  in  $\mathcal{N}$  **do**
  - 5:    $Csum = 0$
  - 6:   **for** each  $k$  in  $\mathcal{K}$  **do**
  - 7:     Calculate  $p_{k,n}^t$  based on (25)
  - 8:     Calculate  $\mathcal{U}_{k,n}^t$  based on (26)
  - 9:     Push  $(k, \mathcal{U}_{k,n}^t)$  into  $h_n$
  - 10:   **end for**
  - 11:   **while**  $h_n$  not empty **do**
  - 12:      $h_n.pop(l, \mathcal{U}_{l,n}^t)$
  - 13:      $Csum = Csum + d_l$
  - 14:     **if**  $Csum \leq D_n$  **then**
  - 15:        $\delta_{l,n}^t = 1$
  - 16:     **end if**
  - 17:   **end while**
  - 18: **end for**
  - 19: Return  $\delta^t$
- 

**Remark 1.** GAPO is actually executable in parallel to a great extent after each RSU  $n$  obtains the required information from the VEC network. In GAPO, the time taken for each RSU  $n$  mainly includes the following parts. The first part is the time taken for RSUs to communicate with each other such that the information required for decision-making can be gathered. The second part is the time taken to sort the services in  $h_n$  and the process can be accomplished with the time complexity  $O(K \log K)$ , where  $K$  is the number of services. The third part is the time taken to traverse  $h_n$  with the aim to find appropriate services to place, and the process can be accomplished with the worst time complexity  $O(K \log K)$ . Note that the first part of the time is usually negligible compared to the other two parts. To sum up, the time complexity of GAPO is  $O(NK \log K)$  without considering the parallel execution of each RSU. The time complexity is much less than the evolutionary algorithms such as Genetic Algorithm (GA).

**Table 1.** Parameter settings.

Para.	Description	Value
$M$	The number of regions	6
$N$	The number of RSUs	20
$K$	The number of services	[25, 65]
$\lambda_{k,m}^t$	Arrival rate for service $k$ in region $m$	[10, 20]
$\mathcal{E}_n$	Energy consumption constraint for RSU $n$	[1.2, 2.2]
$f_n$	The processing frequency	[0.005, 1]GHz
$d_k$	The task-input data size for service $k$	[1, 6]
$s_k$	Average required workload for service $k$	[20, 25]
$D_n$	The caching capability of RSU $n$	[30, 50]
$\varsigma\eta$	The coefficients for energy consumption	1e-12
$\alpha_n$	The static energy consumption	[0, 1)

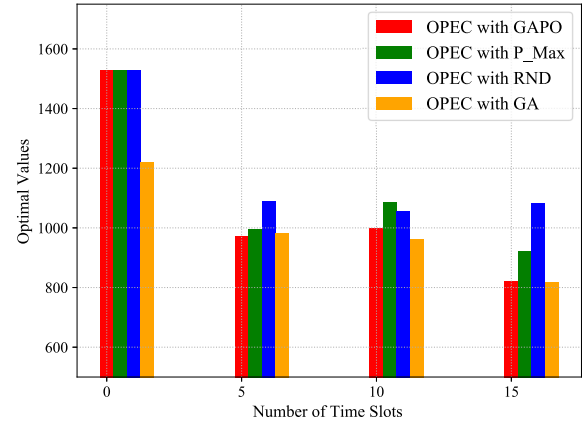
## VI. SIMULATION EVALUATION

### 6.1 Parameter Settings

The main parameters with the default values in the simulation are listed in Table 1. Note that the parameters are set empirically. For example, we assume that the topology of the vehicular edge computing network is built in a random way. In particular, the set of RSUs that can cover the region  $m$  (i.e.,  $\mathcal{R}_m$ ) is determined randomly. Additionally, we assume that the transmission delay can be predicted when task data are transmitted over the backbone network and thus the response delay for task execution in the cloud center can be obtained. For the utility function of services at RSUs, the weights towards the four factors are set to 0.3, 0.2, 0.2, and 0.3, respectively.

### 6.2 Results Analysis and Discussion

Figure 2 shows the performance comparison in terms of the optimal values. In addition to the proposed heuristic, there are actually several approaches for solving the problem  $\mathcal{P}2$ . In particular, we compare our approach with three other approaches in the simulation. The three approaches are the greedy approach, the random approach and the genetic approach, respectively. We denote them by P\_Max, RND, and GA in the simulation, respectively. The greedy approach only adopts the service popularity as the selection criteria. Within each time slot, the services are sorted in descending order based on the defined service popularity. Then, the service with the largest value of popularity is always placed first. The procedure repeats until the violation of constraint conditions, such as (1) happens. The random approach is the simplest way



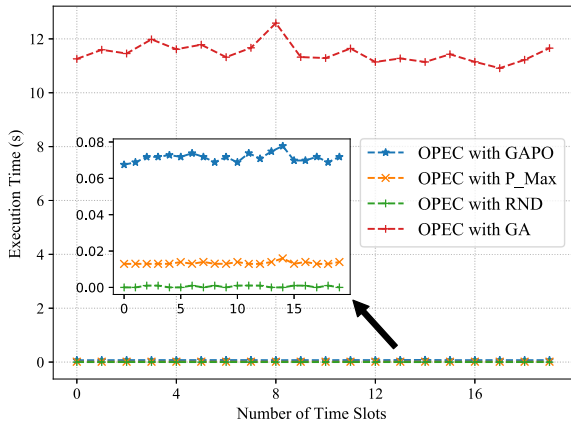
**Figure 2.** The performance comparison with different approaches for  $\mathcal{P}2$ .

to place the services for RSUs within each time slot. During each time slot, services are randomly selected to place as long as no violation happens. The random approach does not consider the feature of the optimization objective or the feature of service requests. The genetic algorithm is especially suitable for binary MKP, owing to its simple deployment and easy implementation. However, the solution accuracy of this approach mainly depends upon the number of iterations, which often turns out to be a time-consuming process. Therefore, it can hardly satisfy the stringent latency requirement in this paper.

From Figure 2, we can observe that the average performance of OPEC with GA is the best among the four approaches, and OPEC with GAPO comes second. In the parameter settings for GA, the crossover and mutation probability are respectively set to 0.2 and 0.02. The population size and the maximal number of iterations are 20 and 30, respectively. In the simulation, some information is the same for all the time slots, which includes  $M$ ,  $N$ ,  $K$ ,  $\mathcal{R}_m$ ,  $\mathcal{E}_n$  and so on. To reflect the dynamics of the vehicular edge computing networks, some information is totally independent in each time slot. Such information is exemplified by  $\lambda_{k,m}^t$ ,  $d_k$  and  $s_k$ .

As a result, we can easily observe that the performance of all approaches wildly fluctuates in different time slots. Even so, our approach is average much better than the random approach and the greedy approach. The random approach can occasionally yield better results, e.g., when the number of time slots is 10. In contrast, the greedy approach, which only considers the

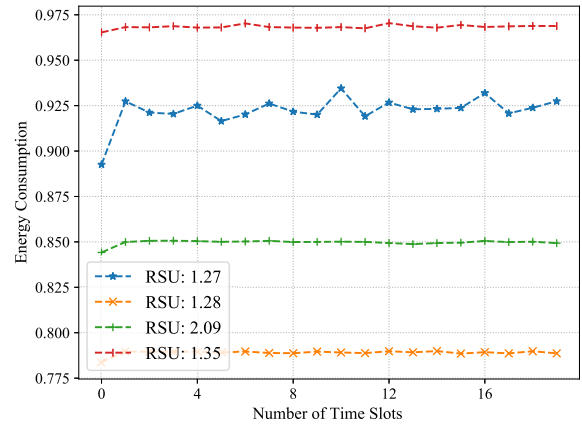




**Figure 3.** The execution time comparison with different approaches for  $\mathcal{P}2$ .

service popularity, never yields the best result but occasionally yields the worst result as denoted in the figure. This is because it neglects two important factors such as  $d_k$  and  $s_k$ , which can seriously affect the value of the optimization objective. For example, based on (9), it is obvious that the service with a larger value of  $s_k$  renders more energy consumption if it is placed at RSU. However, energy consumption has been incorporated into the optimization of the objection function as shown in problem  $\mathcal{P}2$ . Therefore, minimization of the optimization objective is supposed to consider the influence of the averagely required workload  $s_k$ . It shall be noted that the performance of OPEC with GAPO, P\_Max and RND is equally bad compared to OPEC with GA, when the time slot is zero. The service placement profile is generated randomly for the three approaches at the beginning. However, this initial service placement profile can be updated with the increasing iterations in OPEC with GA.

The execution times for different approaches are depicted in Figure 3. Comparatively speaking, the average execution time of OPEC with GAPO is slightly longer than the other two approaches, i.e., OPEC with RND and OPEC with P\_Max. However, these three approaches can almost achieve real-time response in the simulation. Obviously, the execution time of OPEC with GA is much longer than the above three approaches. As shown in Figure 3, this strategy usually takes seconds to get the result in the simulation and thus the execution times of the other three approaches are negligible compared to that of OPEC with GA. Combining with the results in Figure 2, we

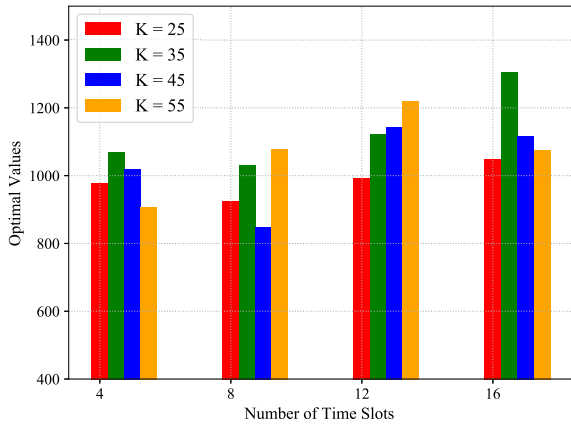


**Figure 4.** The energy consumption along different time slots.

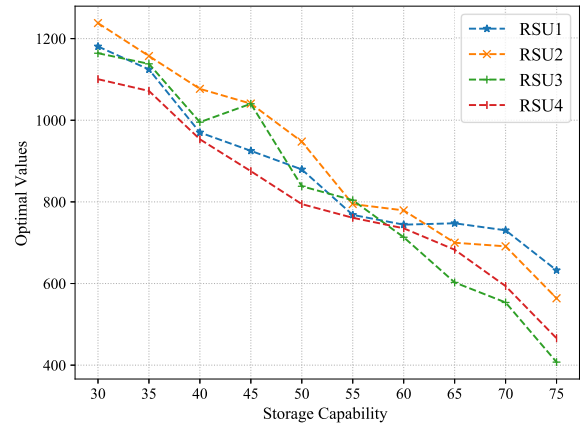
can see that GA helps OPEC achieve the best performance in terms of optimal values at the expense of unbearable time overhead. Considering the strict delay requirement for obtaining the optimal service placement profile, GA-based approach is not appropriate.

Next, we evaluate the average energy consumption at different RSUs along the time slots. The experimental results are shown in Figure 4. On one hand, we need to investigate the energy consumption of each RSU along the time slots; On the other hand, it is more important to check whether the violation happens along the time slots. Specifically, we select 4 from 20 RSUs in the simulation, and the label “RSU: 1.27” means the RSU with its energy threshold being 1.27. First, it is noticeable that the energy consumption of the four RSUs does not exceed their own global energy constraints. As a result, the simulation result has verified the rationality of the theoretical analysis. Second, the energy consumption for RSUs all seem stable along the time slots. Such an observation is acceptable and understandable, since it is in accord with the conclusion drawn in Theorem 2, i.e., the queue backlog is mean rate stable. Last but not least, we notice that the energy consumption for all the RSUs is the least when the time slot is zero. It is explicable since we assume that there are no services to be placed at RSUs in the simulation. Thus, there are no services at RSUs that incur energy consumption at the initial time slot.

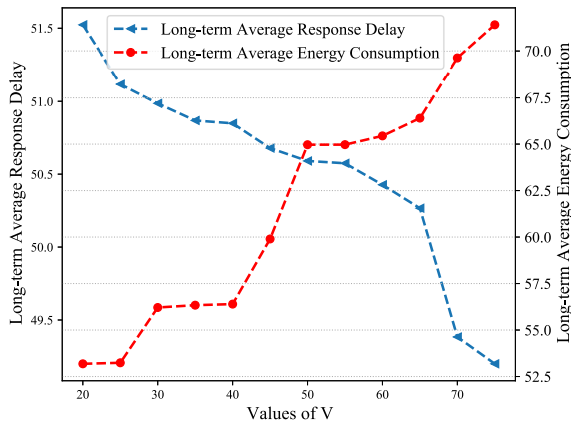
We investigate the influence of the number of services upon the performance of our approach in the next. Figure 5 reveals the optimal values with differ-



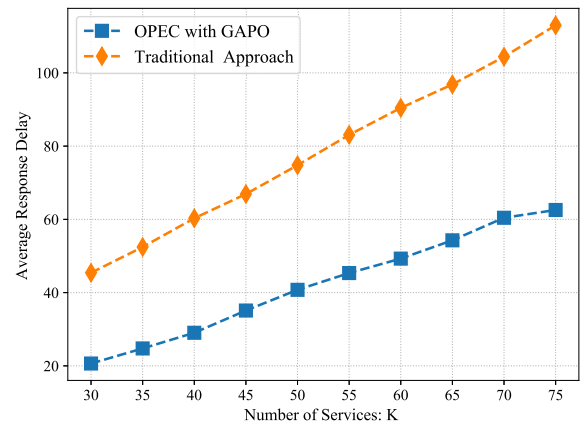
**Figure 5.** The performance comparison under different numbers of services.



**Figure 7.** The performance comparison with different storage capabilities.



**Figure 6.** The performance comparison with different values of  $V$ .



**Figure 8.** The response delays with different numbers of services.

ent numbers of services. Several conclusions can be drawn from the observation as follows. First, it seems that the optimal value of the optimization objective has little to do with the number of services, since there are no deterministic relationships between them. For example, when the number of services is 25, the optimal value in most cases is better than that with the number of services equal to 35 and 45. On the other hand, when the number of services is 55, the optimal value in most cases is also better than that with the number of services equal to 35 and 45, respectively. Accordingly, the number of services is not the more the better, or the less the better, in terms of the optimal values.

The control parameter  $V$  can be used for making a tradeoff between the response latency and energy consumption as mentioned earlier. We study the influence of  $V$  upon the long-term average response de-

lay and the long-term average energy consumption in the simulation. The results are shown in Figure 6. Obviously, the long-term average energy consumption for all RSUs increases with the increasing value of  $V$ . In contrast, the long-term average response delay decreases as  $V$  increases. According to the theoretical analysis, the placement decision profile obtained by the proposed algorithm can gradually approach the best decision profile by increasing the value of  $V$ . However, it is inadvisable to blindly increase the value of  $V$ , considering the ever-increasing average energy consumption for all the RSUs in the long run. It is better to seek a tradeoff between the response delay and energy consumption when  $V$  is set in the simulation.

We investigate the influence of storage capabilities on the performance of OPEC. The simulation results are shown in Figure 7. Obviously, greater storage

capabilities enable more services to be placed at the same time. However, the more services to be placed, the more energy consumption at RSUs. Thus, it is worthwhile investigating the relationships between the performance of the approach and the storage capacity. The simulation result in the figure shows that the performance of OPEC is constrained by the storage capacity at RSUs. For instance, greater storage capacity yields better performance in terms of optimal values.

Since our goal in this paper is to reduce the average response latency for all the requested services by placing services at the edge in advance, we have conducted the last experiment to evaluate our approach compared to the traditional approach that does not apply a service placement strategy. As shown in Figure 8, it is obvious that a better result is yielded when the service placement strategy is applied. In addition, the response delays increase for both approaches as the number of services increases. More importantly, our approach is much better than the traditional approach no matter how the number of services varies.

## VII. CONCLUSION

VEC can cater to the strict delay requirement of vehicular tasks/services, since the computing resources are deployed at the network edge such as RSUs. However, the explosive growth in vehicular offloading requests has caused tremendous pressure on both the edge server and the front-haul links. In this paper, we attempt to optimize the VEC systems by service placement strategy. Specifically, we minimize the average response latency for the requested services along the slotted timeline, while considering multiple constraints such as energy consumption and storage capacity. The greedy heuristic has been incorporated into the drift-plus-penalty-based algorithm for searching for the approximate solution. Extensive simulation has proven that our approach can achieve better performance compared with other approaches in terms of optimal values.

## ACKNOWLEDGMENT

This work was supported by National Natural Science Foundation of China (No. 62071327) and Tianjin Science and Technology Planning Project (No. 22ZYYJC00020).

## References

- [1] X. Li, X. Zhang, *et al.*, "Asynchronous online service placement and task offloading for mobile edge computing," in *18th Annual IEEE International Conference on Sensing, Communication, and Networking, SECON 2021, Rome, Italy, July 6–9, 2021*, pp. 1–9. IEEE, 2021.
- [2] J. Xu, L. Chen, *et al.*, "Joint service caching and task offloading for mobile edge computing in dense networks," in *2018 IEEE Conference on Computer Communications, INFOCOM 2018, Honolulu, HI, USA, April 16–19, 2018*, pp. 207–215. IEEE, 2018.
- [3] C. Tang, X. Wei, *et al.*, "Mobile vehicles as fog nodes for latency optimization in smart cities," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 9, pp. 9364–9375, 2020.
- [4] L. Chen, C. Shen, *et al.*, "Collaborative service placement for edge computing in dense small cell networks," *IEEE Transactions on Mobile Computing*, vol. 20, no. 2, pp. 377–390, 2021.
- [5] C. Tang, C. Zhu, *et al.*, "Toward response time minimization considering energy consumption in caching-assisted vehicular edge computing," *IEEE Internet of Things Journal*, vol. 9, no. 7, pp. 5051–5064, 2021.
- [6] C. Tang, C. Zhu, *et al.*, "Caching assisted correlated task offloading for iot devices in mobile edge computing," in *2021 IEEE Global Communications Conference (GLOBECOM)*, pp. 1–6. IEEE, 2021.
- [7] R. Xie, Q. Tang, *et al.*, "When serverless computing meets edge computing: Architecture, challenges, and open issues," *IEEE Wireless Communications*, vol. 28, no. 5, pp. 126–133, 2021.
- [8] H. Hou, H. Jin, *et al.*, "Cost efficient edge service placement for crowdsensing via bus passengers," *Mobile Networks and Applications*, vol. 26, pp. 899–908, 2021.
- [9] M. Ayoubi, M. Ramezanzpour, *et al.*, "An autonomous iot service placement methodology in fog computing," *Softw. Pract. Exp.*, vol. 51, no. 5, pp. 1097–1120, 2021.
- [10] T. Huang, W. Lin, *et al.*, "An ant colony optimization-based multiobjective service replicas placement strategy for fog computing," *IEEE Transactions on Cybernetics*, vol. 51, no. 11, pp. 5595–5608, 2021.
- [11] X. Xu, B. Shen, *et al.*, "Edge server quantification and placement for offloading social media services in industrial cognitive iot," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 4, pp. 2910–2918, 2021.
- [12] Y. Hao, M. Chen, *et al.*, "Deep reinforcement learning for edge service placement in softwarized industrial cyber-physical system," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 8, pp. 5552–5561, 2021.
- [13] Y. Chen, S. Zhang, *et al.*, "Locus: User-perceived delay-aware service placement and user allocation in mec environment," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 7, pp. 1581–1592, 2021.
- [14] J. Ortín, J. R. Gállego, *et al.*, "On optimizing network function placement for multicast group call service provision in lte iops networks," *IEEE Access*, vol. 9, pp. 160 897–160 916, 2021.
- [15] B. Yuan, S. Guo, *et al.*, "Joint service placement and request routing in mobile edge computing," *Ad Hoc Networks*, vol. 120, p. 102543, 2021.
- [16] J. Panadero, M. Selimi, *et al.*, "A two-stage multi-criteria

optimization method for service placement in decentralized edge micro-clouds,” *Future Generation Computer Systems*, vol. 121, pp. 90–105, 2021.

- [17] P. Han, Y. Liu, *et al.*, “Interference-aware online multicomponent service placement in edge cloud networks and its ai application,” *IEEE Internet of Things Journal*, vol. 8, no. 13, pp. 10 557–10 572, 2021.
- [18] V. Farhadi, F. Mehmeti, *et al.*, “Service placement and request scheduling for data-intensive applications in edge clouds,” *IEEE/ACM Transactions on Networking*, vol. 29, no. 2, pp. 779–792, 2021.
- [19] Z. Ning, P. Dong, *et al.*, “Distributed and dynamic service placement in pervasive edge computing networks,” *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 6, pp. 1277–1292, 2021.
- [20] X. He, R. Jin, *et al.*, “Joint service placement and resource allocation for multi-uav collaborative edge computing,” in *2021 IEEE Wireless Communications and Networking Conference (WCNC)*, pp. 1–6. IEEE, 2021.
- [21] Z. Liu, C. Zhan, *et al.*, “Robust edge computing in uav systems via scalable computing and cooperative computing,” *IEEE Wireless Communications*, vol. 28, no. 5, pp. 36–42, 2021.
- [22] C. Tang, C. Zhu, *et al.*, “Integration of UAV and fog-enabled vehicle: Application in post-disaster relief,” in *25th IEEE International Conference on Parallel and Distributed Systems, ICPADS 2019, December 4-6, 2019*, pp. 548–555. IEEE, 2019.
- [23] L. Georgiadis, M. J. Neely, *et al.*, “Resource allocation and cross-layer control in wireless networks,” *Foundations & Trends in Networking*, vol. 1, no. 1, pp. 1–144, 2006.
- [24] Neely and J. Michael, “Stochastic network optimization with application to communication and queueing systems,” *Synthesis Lectures on Communication Networks*, vol. 3, no. 1, p. 211, 2010.
- [25] K. Bernhard and J. Vygen, “Combinatorial optimization: Theory and algorithms,” *Springer, Third Edition, 2005*, 2008.

## Biographies



**Chaogang Tang** received his B.S. degree from the Nanjing University of Aeronautics and Astronautics, Nanjing, China, and Ph.D. degree from the School of Information Science and Technology, University of Science and Technology of China, Hefei, China, and the Department of Computer Science, City University of Hong Kong, Hong Kong, China, under a joint Ph.D. program, in 2012. He is now with the China University of Mining and Technology. His research interests include mobile cloud computing, fog computing, Internet of Things, and big data.



**Yubin Zhao** received his B.S. and M.S. in 2007 and 2010 respectively from Beijing University of Posts and Telecommunications (BUPT), Beijing, China. He received his Ph.D. degree in computer science in 2014 from Freie Universität Berlin (FU Berlin), Berlin, Germany. He is currently an Associate Professor with the School of Microelectronics Science and Technology, Sun Yat-Sen University, Zhuhai, China. His current research interest includes wireless power transfer, indoor localization and target tracking.



**Huaming Wu** received the B.E. and M.S. degrees from Harbin Institute of Technology, China in 2009 and 2011, respectively, both in electrical engineering. He received the Ph.D. degree with the highest honor in computer science at Freie Universität Berlin, Germany in 2015. He is currently an associate professor at the Center for Applied Mathematics, Tianjin University, China. His research interests include mobile cloud computing, edge computing, Internet of Things, deep learning, complex network, and DNA storage.