

Towards Response Time Minimization Considering Energy Consumption in Caching Assisted Vehicular Edge Computing

Chaogang Tang, *Member, IEEE*, Chunsheng Zhu, *Member, IEEE*, Huaming Wu, *Member, IEEE*, Qing Li, *Senior Member, IEEE*, and Joel J. P. C. Rodrigues, *Fellow, IEEE*

Abstract—The advent of vehicular edge computing (VEC) has generated enormous attention in recent years. It pushes the computational resources in close proximity to the data sources and thus caters for the explosive growth of vehicular applications. Owing to the high mobility of vehicles, these applications are of latency-sensitive requirements in most cases. Accordingly, such requirements still pose a great challenge to the computing capabilities of VEC, when these applications are outsourced and executed in VEC. Against this backdrop, we propose a new mathematical model which respectively generalizes the computation and communication models, and applies application oriented caching into VEC in this paper. Based on this model, a new strategy is further proposed to optimize the average response time of applications over an infinite time-slotted horizon for VEC. A long-term energy consumption constraint is imposed to guarantee the stability of VEC system, and the Lyapunov optimization technology is adopted to tackle this constraint issue. Two greedy heuristics are put forward to help find the approximate optimal solution in the drift-plus-penalty based algorithm. Extensive experiments have been conducted to evaluate the response time and energy consumption in the caching assisted VEC. The simulation results have shown that the proposed strategy can dramatically optimize the average response time while satisfying the long-term energy consumption constraint.

Index Terms—Vehicular edge computing, caching, service provisioning, Lyapunov optimization, greedy heuristics.

Copyright (c) 20xx IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org. This work is partially supported by the National Key R&D Program of China (2020YFB2104301), by Chinese National Research Fund (NSFC) Key Project No. 61532013 and No.61872239, the project “Network Communication Intelligent Core Chip Design and Core Software (PCL2021-A08)”, and the Project “Beihang Beidou Technological Achievements Transformation and Industrialization Funds (BARI2005)”. This work is also partially supported by FCT/MCTES through national funds and when applicable co-funded EU funds under the Project UIDB/50008/2020, and by Brazilian National Council for Research and Development (CNPq) via Grant No. 313036/2020-9. (Corresponding author: Chunsheng Zhu)

C. Tang is with the School of Computer Science and Technology, China University of Mining and Technology, 221116, Xuzhou, China. C. Tang is also with Mine Digitization Engineering Research Center of the Ministry of Education, China University of Mining and Technology, Xuzhou 221116, China (e-mail: cgtang@cumt.edu.cn).

C. Zhu is with the SUSTech Institute of Future Networks, Southern University of Science and Technology, 518055, Shenzhen, China. C. Zhu is also with the PCL Research Center of Networks and Communications, Peng Cheng Laboratory, 518055, Shenzhen, China (e-mail: chunsheng.tom.zhu@gmail.com).

H. Wu is with the Center for Applied Mathematics, Tianjin University, 300072, Tianjin, China (e-mail: whming@tju.edu.cn)

Q. Li is with the Hong Kong Polytechnic University, Hong Kong, China (e-mail: qing-prof.li@polyu.edu.hk)

J. Rodrigues is with Federal University of Piauí, Teresina - Pi, Brazil. J. Rodrigues is also with Instituto de Telecomunicações, Portugal (joelljr@ieee.org).

Manuscript received ; revised .

I. INTRODUCTION

THE rapid development of information and communication technology gives rise to the prosperities of smart vehicles [1], [2]. In addition to the ability to communicate with each other, smart vehicles are also equipped with computer facilities to perform vehicular applications. However, a large number of vehicular applications have created enormous pressure on such “computers with the wheels”. On the other hand, the advent of vehicular edge computing (VEC) [3] has generated enormous attention recently. In comparison with the limited capabilities of the on-board computers, there are more computational resources in VEC. Furthermore, VEC pushes these resources in close proximity to the data sources. For instance, the edge servers in VEC are usually deployed at the logical edge of networks, e.g., road side unit (RSU). Vehicles can outsource their applications to RSU via Vehicle-to-Infrastructure (V2I) communication technologies. Hence, this computing paradigm caters for the increasingly complicated requirements of vehicular applications.

Additionally, VEC can alleviate traffic congestion in the network core and lays profound foundation for the development of vehicular applications in smart transportation. Despite the benefits brought for vehicular applications, VEC is still confronted with a big challenge. Most of the vehicular applications are time-sensitive such as route navigating [4], owing to the high mobility of vehicles. However, it is very challenging for VEC to satisfy the strict latency requirements, due to the fact that the computational resources in VEC are not inexhaustible compared to cloud computing. The response latency reduction comes at the expense of the computational resource constraint in VEC.

Therefore, significant efforts are still needed to improve the performance of VEC systems. For this purpose, we propose to combine the application/service oriented caching with VEC in this paper. By caching the most frequently outsourced applications in VEC, the response time can be tremendously reduced. In our view, caching assisted VEC can greatly improve the efficiency of application outsourcing and task performing, and thus fulfill the key requirements of the computationally intensive and time-sensitive applications. However, different from content delivery in caching enabled information centric networks [5], [6], application outsourcing in the caching assisted VEC is much more complicated in terms of response time and energy consumption [7]–[10].

For example, virtual machine (VM) and lightweight container Docker [11] are the two most widely used frameworks in the provision of application outsourcing at the edge. Both of them will temporarily initialize virtual environments for the offloaded applications. If the applications are cached at the server, it will be better to maintain the running state of virtual environments for a certain period of time, with the aim to avoid additional time overheads on virtual initialization. This way, however, will incur additional energy consumptions. Since the energy consumption can serve as an important performance indicator to evaluate the VEC systems, the tradeoff between the energy consumption and application caching strategies should be carefully considered. On another hand, the evaluation of the caching assisted VEC system, such as the response time optimization, usually needs a long-term process. Therefore, it becomes significant and yet very challenging to design an appropriate application/service caching strategy in VEC.

In this paper, we focus on the response time minimization while considering the energy consumption for application outsourcing in the caching assisted VEC. In particular, we list the major contributions of this paper, as below.

- We propose a generic approach to improve the performance of application outsourcing in the caching assisted VEC. Specifically, we mathematically formulate the optimization problem, with the aim to minimize the average response time of the outsourced applications over a long time-slotted horizon in this paper.
- Owing to the difficulty of solving the optimization problem while satisfying the energy consumption over an infinite time-slotted horizon, we introduce the Lyapunov optimization technology in this paper. By doing so, the long-term energy constraint can be converted into per-slot ones. Furthermore, a dynamic algorithm based on *drift-plus-penalty* term is proposed to obtain the approximate optimal solution over the accumulated time slots.
- Extensive experiments are carried out to evaluate the response time and energy consumption of the proposed caching strategy in VEC. Simulation results show that the proposed application outsourcing strategy in the caching assisted VEC can achieve better performance while satisfying the time slot spanned energy constraint.

The rest of paper is organized as follows. Representative works are studied and discussed in Section II. In Section III, we apply application-oriented caching into VEC and formally establish a mathematical model to minimize the response time while considering the energy consumption in the caching assisted VEC. In Section IV, we apply the Lyapunov optimization technology to tackle the time slot spanned energy constraint, and further put forward a drift-plus-penalty based algorithm which leverages two greedy heuristics to find the approximate optimal solution in the caching assisted VEC. The simulation results are reported and discussed in Section V, followed by a conclusion in Section VI.

II. RELATED WORKS

With the advent of smart transportation and its corresponding sub-ecosystems such as smart vehicle and RSU, the

vehicle-loaded computers have gained rapid development. For example, the computational resources can be extended to such entities, enabling various applications and services running at the vehicles and RSUs, while satisfying the strict response time requirement. VEC has been considered as a promising and efficient approach to relieve the burden of backbone networks [12], because a vast amount of data generated by numerous smart terminals can be processed locally or at the logical edge of networks, with no need for task offloading to the remote cloud center via the core networks.

Recently, the task offloading and service outsourcing in VEC have indeed attracted increasing attention by virtue of its advantages [13]–[15]. We plan to investigate the existing literature on VEC from two aspects in what follows.

A. Caching Based VEC

Inspired by the content-centric mobile edge caching, researchers turn their attention to the cache enabled task offloading in VEC [7], [8], with an purpose of further optimizing energy consumption, response time, and so on. For instance, if the tasks are cached in VEC, vehicles do not need to offload the tasks any longer and hence the transmission delay and execution time in the edge can be omitted [9], [10].

Authors in [7] suggest that a small number of services can be cached in size-limited edge server so as to improve QoS with regards to (w.r.t.) response time. To investigate the performance of the caching strategies, the optimization problem is formulated to minimize the energy consumption in the long run, i.e., across different time slots. An efficient online algorithm is proposed to tackle these issues in mobile edge computing systems, which tries to jointly optimize the dynamic service caching and task offloading.

The computation-intensive and time-sensitive tasks are becoming dominating recently, which benefits from increasingly popular applications such as augment reality and in-car gaming. Authors in [8] consider the task caching as a promising solution to satisfy the low latency requirement. In this regard, they propose to cache the completed task application and corresponding data in the edge-cloud, and try to jointly optimize the task caching and offloading with multiple constraints. An alternating iterative algorithm is proposed to tackle this problem and simulation results prove that it indeed outstands other approaches.

To realize efficient information dissemination, authors in [16] construct a hierarchical end-edge framework which combines data delivery, computation offloading and content caching. The network overhead is selected as the performance indicator to evaluate this framework and the related strategy. In particular, they model the optimization problem as a mixed integer non-linear programming problem and solve it based on the deep deterministic policy gradient. Simulation results have revealed its advantages compared to other benchmark approaches.

A mass of vehicles can serve as the intermediate nodes to provide caching services for the requests in proximity. However, several issues should be addressed, which include the QoS, incentives and privacy. To this end, authors in [17] combine the

deep reinforcement learning with the permissioned blockchain to achieve an intelligent and secure content caching in the VEC system and network. Simulation results have indicated that the approach is much better than two other benchmark methods.

B. Response Time and Energy Consumption about VEC

Vehicles can serve as the computing nodes to provide the computational resources to cater for the time-sensitive applications in the vicinity [8], [18]–[20]. For example, with the increasing number of applications, various computational entities are welcome with the aim to mitigate the high demands of computing resources in smart city. In [18], authors first proposed that vehicles with idle computing resources should act as the infrastructures to perform multiple functionalities such that vehicles can play more important roles in communication and computation. Specifically, they have discussed several kinds of communication and computational infrastructures where either moving or parked vehicles can serve as the computing entities for response latency optimization or improving the throughput of the computing paradigm.

Similar works can also be found in [19] where authors suggest that the ride-share taxis can serve as the infrastructure to provide both communication and computation resources for the time-sensitive applications. To be specific, terminal users can leverage these taxis to support live video streaming, e.g., taxis can proactively receive video chunks. Furthermore, they formulate the problem as a coverage optimization problem and solve it by the strategy that is applied to the set cover problem (SCP). Simulation results have proven its advantages compared to other strategies.

On the other hand, the resource-hungry and time-sensitive in-car applications have put huge pressure on the limited computing capabilities of vehicles. Meanwhile, VEC has shown great potentials to tackle such an issue, e.g., by outsourcing the vehicular applications to the ubiquitous vehicular edge servers such as RSU [14], [20], [21].

A multi-user multi-server VEC system is considered in [21] where the load balance-aware resource allocation is studied when tasks are offloaded. Authors have regarded this problem as a mixed integer nonlinear programming problem for the system utility optimization. In particular, they divide it into two subproblems and design efficient algorithms for them.

In [22], authors strive to establish an offloading system in VEC with the help of deep reinforcement learning technology. Specifically, the finite Markov chains are applied to modelling the communication and computation states. The optimization problem is modeled mathematically in the hope to improve the Quality of Experience (QoE) of users. Similar to [21], the problem is also decoupled into two sub-optimization problems, and solved by a low-complexity algorithm.

Owing to the resource-hungry feature of the tasks, a certain portion of tasks pertaining to the strict latency requirement are encouraged to offload to the external entities for execution. VEC can play an important role in assisting such kinds of task execution. To predict the resource consumption in the edge nodes more accurately, authors in [23] try to achieve an adaptive selection of machine learning algorithms, by

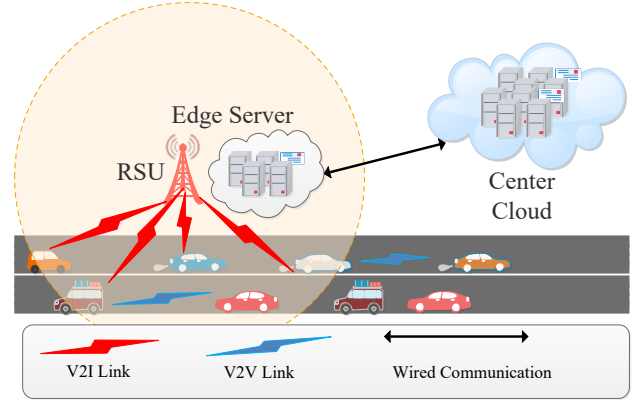


Fig. 1. Application Scenario

using a two-stage meta-learning based approach and extracting meta-features on database. Multiple offloading decisions are considered in [15] where tasks are offered three options for the execution, i.e., local host, cloudlet and cloud center. To achieve better performance of the cloudlet-cloud offloading, a mathematical problem is proposed that aims to optimize the energy consumption in the long run.

For the same task with different task-input parameters, the computational result could be of great difference. Accordingly, the cached computational result may not benefit the following requestors even if the tasks are the same. Few of the aforementioned works have considered this case. Compared to these works, we respectively generalize the computation, communication and caching model to accommodate this case, and aim to minimize the response time with the consideration of the energy consumption constraint in the caching assisted VEC.

III. PROPOSED CACHING ASSISTED VEC MODEL

We consider a VEC scenario consisting of one edge server S and multiple moving vehicles, as shown in Fig. 1. The edge server can be deployed at a road side unit (RSU) which directly connects to the cloud server on one hand, and communicates with nearby vehicles using the vehicle-to-infrastructure (V2I) communication technology on the other hand [24]. Due to the limited computing capabilities, vehicles can outsource the in-car applications to S for execution, with the aim to alleviate the high demand of computing resources. We assume that the offloaded applications belong to a set of K applications, indexed by $\mathcal{A} = \{a_1, a_2, \dots, a_K\}$, and each application a_k can be described as $a_k = (d_k, s_k)$ where d_k denotes the average task-input size (e.g., the processing codes and user related parameters) which need to be offloaded over the wireless channels, and s_k the average number of CPU cycles for accomplishing a_k .

We consider a caching assisted service outsourcing in VEC in this paper, and expect that the response time can be drastically reduced by caching those most frequently requested services/applications at the edge server. For instance, when an application a_k that a vehicle wants to offload is already cached in S , the vehicle does not need to transmit d_k any

longer, but to wait for the result to be returned from S . It usually takes a long-term process to evaluate the overall performance of the caching enabled VEC system. To this end, time is divided into a set of discrete time slots, indexed by $\{0, 1, \dots, T-1\}$ and each time slot has a duration π . The service caching strategy can be updated within each time slot. Note that we use the three words “applications”, “services” and “tasks” exchangeably hereafter.

Owing to the stochastic nature of vehicles, we assume that the arrival of the application offloading requests for a_k in time slot t follows a Poisson process with arrival rate λ_k^t [25]. We define $\delta_k^t \in \{0, 1\}$ as a binary variable to denote whether a_k is cached at S in the t th time slot. If a_k is cached at S , $\delta_k^t = 1$, and 0, otherwise. Thus, the caching decision profile in time slot t can be expressed as $\delta^t = (\delta_1^t, \delta_2^t, \dots, \delta_K^t)$. And we have $\sum_{k=1}^K \delta_k^t d_k \leq \mathcal{C}$, $\forall t \in [0, T-1]$ where \mathcal{C} denotes the cache size of S .

Specifically, Fig. 2 shows the sketch of caching assisted service outsourcing in VEC. For the vehicles with the offloading requests in the vicinity, they can disseminate such requests to VEC (RSU) together with the beacon information. RSU gathers the offloading requests and decides which services should be cached according to the evaluation metric (illustrated later). The caching profile determined by RSU is then sent back to the vehicles. Based on the caching decision, vehicles will perform different offloading operations. For example, if the requested service has been cached in VEC, the vehicle only needs to offload the context information (e.g., the user related parameters). In this case, the time taken to offload the context information is negligible compared to the task itself. On the other hand, if the requested service is not cached in VEC, the vehicle will offload the task to RSU for the execution. Our work in this paper focuses on the caching decision making after the information on the requested services is gathered. It shall be noted that at the beginning of each time slot, the caching decision profile can be updated based on the beacon information disseminated from vehicles. For convenient discussion, some key notations to be used are shown in Table I.

A. Service Caching Model

Intuitively, the more the number of services which are cached, the more the response time can be reduced. However, due to the limited cache size of S , it is impractical to cache all the services at the same time. On another hand, from the viewpoint of response latency reduction, it is better for the edge server to maintain the virtual environment for a while after a service is cached. However, the corresponding energy consumptions at S will substantially increase as the number of cached services increases. Therefore, we should judiciously determine which services to be cached at S , for the purpose of response latency reduction while considering the energy consumption.

If application a_k is not cached at S in time slot t , vehicles should offload it to S and wait for the execution result. In this case, the average response time $l_{k,rl}^{t,no}$ consists of the transmission response $l_{k,trs}^{t,no}$ and execution response $l_{k,exe}^{t,no}$,

TABLE I
NOTATIONS

Notation	Description
K	The number of vehicular applications
S	The edge server deployed at RSU
\mathcal{A}	The set of applications to be deployed
d_k	The average size of input data
s_k	The average number of CPU cycles needed for a_k
T	The total number of time slots
π	The duration for each time slot
\mathcal{C}	The cache size of S
β_k^t	The number of channels assigned to a_k
r_k^t	The average transmission rate over the channel
λ_k^t	The arrival rate for a_k in time slot t
f_e^t	The average processing frequency of S in time slot t
κ	The effective switched capacitance coefficient
ε	The number of cycles for one task-input bit performing
γ_k	The static power consumption at S for a_k
E_t^{max}	The energy consumption constraint for time slot t
L_t^{max}	The response time constraint for time slot t
Q	The global energy constraint across different slots

each of which can not be neglected. However, the time taken to return the result from S can be omitted due to the negligible result size [26]. We need to calculate $l_{k,trs}^{t,no}$ and $l_{k,exe}^{t,no}$, respectively. The average transmission delay $d_{k,trs}^{t,no}$ for a_k over the wireless channel in slot t can be expressed as:

$$d_{k,trs}^{t,no} = \frac{d_k}{\beta_k^t r_k^t} \quad (1)$$

where β_k^t denotes the number of channels assigned to a_k and r_k^t the average transmission rate over the channel. To simplify the discussion, we assume that r_k^t is a constant which can be obtained based on the historical statistics. The arrival of tasks at the transmission channel also follows a Poisson process based on the aforementioned assumption. The arrival rate at the transmission queue is λ_k^t and the service rate is $\pi \beta_k^t r_k^t / d_k$. Based on the $M/M/1$ queueing model, the average queuing delay in the transmission channel can be calculated as [27]:

$$d_{k,q}^{t,no} = \frac{\lambda_k^t d_k^2}{\pi \beta_k^t r_k^t (\pi \beta_k^t r_k^t - \lambda_k^t d_k)} \quad (2)$$

Since the average transmission response consists of the average queuing delay in the transmission channel and the following average transmission delay, we have:

$$l_{k,trs}^{t,no} = d_{k,q}^{t,no} + d_{k,trs}^{t,no} \quad (3)$$

The execution delay of a_k at S can be calculated as

$$d_{k,exe}^{t,no} = d_{k,init}^{t,no} + \frac{s_k}{f_e^t} \quad (4)$$

where f_e^t is the average processing frequency of S in time slot t and $d_{k,init}^{t,no}$ is the time taken to initialize the virtual machine (VM) for a_k in time slot t .

To obtain the queueing delay of a_k , information about other applications $\mathcal{A} \setminus \{a_k\}$ in time slot t is also required. Based on the queueing model, the arrival rate of all the services $\tilde{\lambda}^t = \sum_{k=1}^K \lambda_k^t$ at the task queue of S also follows a Poisson process. Thus, the average execution time (i.e., service time) for a service is $\sum_{k=1}^K \lambda_k^t s_k / \pi f_e^t \tilde{\lambda}^t$. Thus, the service rate for

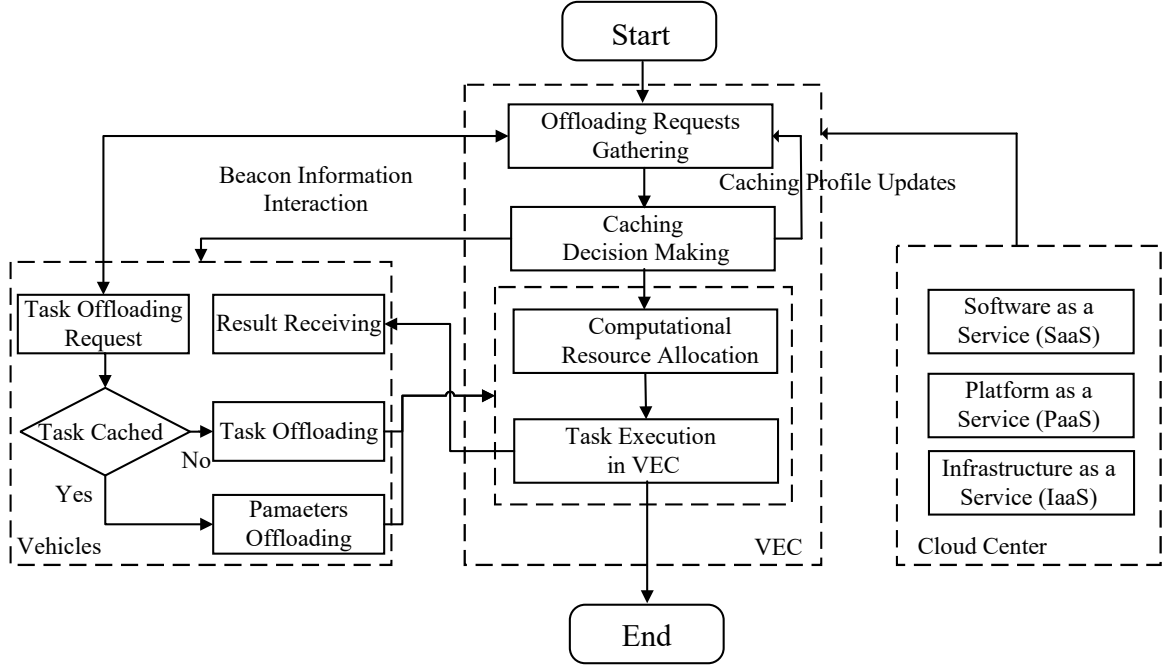


Fig. 2. A sketch of caching assisted service outsourcing in VEC

each VM is $\pi f_e^t \tilde{\lambda}^t / \sum_{k=1}^K \lambda_k^t s_k$. The number of VMs in S is assumed to be m . According to the $M/M/s$ queueing theory, the average queueing delay for a_k is given as:

$$d_{k,aq}^{t,no} = \frac{m\rho}{\tilde{\lambda}^t} + \frac{m^m \rho^{m+1} P_0}{\tilde{\lambda}^t m!} \quad (5)$$

where $\rho = \tilde{\lambda}^t \sum_{k=1}^K \lambda_k^t s_k / (m \pi f_e^t \tilde{\lambda}^t)$ and $P_0 = 1 / (\sum_{k=0}^{m-1} (m\rho)^k / k! + (m\rho)^m / (m!(1 - \rho)))$. Thus, the average execution response time $l_{k,exe}^{t,no}$ can be expressed as:

$$l_{k,exe}^{t,no} = d_{k,aq}^{t,no} + d_{k,exe}^{t,no} \quad (6)$$

Furthermore, the average energy consumption for a_k at S can be given as

$$e_k^{t,no} = \kappa \varepsilon s_k (f_e^t)^2 \quad (7)$$

where κ and ε represent the effective switched capacitance coefficient and the number of cycles required for one task-input bit performing at S , respectively.

If a_k is cached at S in time slot t , it is unnecessary for vehicles to offload the corresponding tasks to S . In this case, the average response time of a_k denoted by $l_{k,rl}^{t,c}$ is directly equal to the execution response time $l_{k,exe}^{t,c}$. In most of the existing works, $l_{k,exe}^{t,c} = 0$ holds, as long as the service is cached. The assumption behind it is that the computational result is applicable to any offloading request irrespective of context information and user related parameters. As a contrast, we allow for the diversity of context information and user parameters. As a result, $l_{k,exe}^{t,c}$ needs to be recalculated even if the corresponding service is cached. The execution delay of a_k in time slot t is: $d_{k,exe}^{t,c} = s_k / f_e^t$. Since the running state of virtual environment is maintained when the service is cached, the initialization time can be omitted, which is

different from $d_{k,exe}^{t,no}$. The average queueing delay for a_k , i.e., $d_{k,aq}^{t,c}$, is the same as $d_{k,aq}^{t,no}$. Thus, when the corresponding service is cached, the average execution response consisting of the average queueing delay and average execution delay can be expressed as

$$l_{k,exe}^{t,c} = d_{k,aq}^{t,c} + d_{k,exe}^{t,c} \quad (8)$$

The corresponding energy consumption can be calculated as

$$e_k^{t,c} = \gamma_k + \kappa \varepsilon s_k (f_e^t)^2 \quad (9)$$

where γ_k is the static power consumption caused by the maintenance of the virtual environment for a_k regardless of the workloads. Therefore, we have the expected energy consumption E_t in time slot t as follows:

$$E_t(\delta^t) = \sum_{k=1}^K (1 - \delta_k^t) e_k^{t,no} + \delta_k^t e_k^{t,c} \quad (10)$$

B. Problem Formulation

Our goal in this paper is to minimize the average response time over a long time-slotted horizon while keeping the energy consumption at S below the given threshold. Accordingly, the optimization problem is formulated as follows:

$$(P1) \quad \min_{\delta^t, \forall t} \quad \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{k=1}^K [(1 - \delta_k^t) l_{k,rl}^{t,no} + \delta_k^t l_{k,rl}^{t,c}] \quad (11)$$

$$s.t. \quad \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} E_t(\delta^t) \leq Q \quad (12)$$

$$\sum_{k=1}^K \delta_k^t d_k \leq C \quad \forall t \in [1, T] \quad (13)$$

$$E_t(\delta^t) \leq E_t^{max} \quad \forall t \in [1, T] \quad (14)$$

$$\sum_{k=1}^K [(1 - \delta_k^t) l_{k,rl}^{t,no} + \delta_k^t l_{k,rl}^{t,c}] \leq L_t^{max} \quad \forall t \in [1, T] \quad (15)$$

where L_t^{max} and E_t^{max} denote the maximum response time and energy consumption allowed for time slot t , respectively. Specifically, conditions (14) and (15) place L_t^{max} and E_t^{max} on the per-slot response time and energy consumptions, respectively. Condition (12) represents that the average energy consumption over the accumulated time slots should not go beyond the energy consumption constraint Q . Constraint (13) ensures that the amount of services to be cached should not exceed the cache size of S .

Challenges. Exhaustive search over the potential solution space is prohibitively costly, since it takes the exponential time to determine the best caching profile for each time slot. Furthermore, to optimally solve problem $P1$ requires the future information (e.g., service distributions in all time slots). Based on these information, the optimal solution can be obtained in an off-line way. However, it is highly difficult to predict such information beforehand in reality. Accordingly, these challenges necessitate an online approach which can efficiently make service caching decisions without the future information.

IV. RESPONSE TIME MINIMIZATION STRATEGY CONSIDERING ENERGY CONSUMPTION IN CACHING ASSISTED VEC

A. Lyapunov-Based Online Caching Decision Algorithm

Directly solving problem $P1$ is challenged by not only the prediction on the future information but also the entire energy consumption constraint across the time slots. Accordingly, the Lyapunov optimization framework is adopted in this paper to tackle such challenges. We construct a dynamic energy migration queue to assist the caching decision making while satisfying the long-term energy constraint. Let $q(0) = 0$, and this queue can be recursively defined as:

$$q(t+1) = \max[q(t) - Q, 0] + E_t(\delta^t) \quad (16)$$

where $q(t)$ is the queue backlog in time slot t that reflects the deviation of current energy consumption from the energy constraint Q . Hence, the larger the value of $q(t)$, the sharper the deviation.

Lemma 1: Given the queue backlog $q(t)$, the following inequality holds

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[E_t(\delta^t) - Q] \leq \frac{1}{T} \mathbb{E}[q(T)]$$

Proof We can briefly prove it as follows. If $q(t) \geq Q$, according to the definition of $q(t)$, we have $q(t+1) = \max[q(t) - Q, 0] + E_t(\delta^t) = q(t) - Q + E_t(\delta^t)$, thus, $E_t(\delta^t) - Q = q(t+1) - q(t)$. If $q(t) < Q$, we have $q(t+1) = \max[q(t) - Q, 0] + E_t(\delta^t) = E_t(\delta^t)$ and $E_t(\delta^t) - Q < q(t+1) - q(t)$. Accordingly, $E_t(\delta^t) - Q \leq q(t+1) - q(t)$

$\forall t \in \{0, \dots, T-1\}$. Take the expectation of this inequality and further summarize it over $\forall t \in \{0, \dots, T-1\}$, we have:

$$\sum_{t=0}^{T-1} \mathbb{E}[E_t(\delta^t) - Q] \leq \sum_{t=0}^{T-1} \mathbb{E}[q(t+1) - q(t)] = \mathbb{E}[q(T)]$$

Thus,

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[E_t(\delta^t) - Q] \leq \frac{1}{T} \mathbb{E}[q(T)]$$

□

The Lyapunov function can be defined as $L(q(t)) = \frac{1}{2}q^2(t)$. To ensure the strong stability of this migration queue, the increment between two consecutive states should be as small as possible. To this end, we define the *Lyapunov drift* as $\Delta(q(t)) \triangleq \mathbb{E}[L(q(t+1)) - L(q(t))|q(t)]$. We can determine the upper bound of $\Delta(q(t))$, based on the Lemma from [28], given below.

Lemma 2: Assume A, B, C and m are non-negative real numbers and $C = \max\{B - m, 0\} + A$, then $C^2 \leq A^2 + B^2 + m^2 - 2B(m - A)$.

Based on this Lemma, we have

$$\begin{aligned} \Delta(q(t)) &= \mathbb{E}[L(q(t+1)) - L(q(t))|q(t)] \\ &= \frac{1}{2} \mathbb{E}[(\max[q(t) - Q, 0] + E_t(\delta^t))^2 - q^2(t)|q(t)] \\ &\leq \frac{1}{2} \mathbb{E}[Q^2 + E_t^2(\delta^t) + 2q(t)(E_t(\delta^t) - Q)|q(t)] \\ &= \frac{1}{2} Q^2 + \mathbb{E}\left[\frac{E_t^2(\delta^t)}{2} - Qq(t) + q(t)E_t(\delta^t)|q(t)\right] \\ &= D - Qq(t) + \mathbb{E}[q(t)E_t(\delta^t)|q(t)] \end{aligned} \quad (17)$$

where $D = \frac{Q^2}{2} + \mathbb{E}\left[\frac{E_t^2(\delta^t)}{2}|q(t)\right] \leq \frac{Q^2}{2} + \mathbb{E}\left[\frac{(E_t^{max})^2}{2}|q(t)\right] = \frac{Q^2}{2} + \frac{(E_t^{max})^2}{2} \triangleq B$. As a result, $\Delta(q(t)) \leq B - Qq(t) + \mathbb{E}[q(t)E_t(\delta^t)|q(t)]$. It is easily observed that the upper bound of $\Delta(q(t))$ does not need the future information such as $q(t+1)$ at time slot t . Therefore, we can convert the time slot spanned energy constraint into per-slot ones by means of the Lyapunov optimization technology. Specifically, the *drift-plus-penalty* term in each time slot is given as:

$$\begin{aligned} &\Delta(q(t)) + V\mathbb{E}[L_t(\delta^t)|q(t)] \\ &\leq B - Qq(t) + \mathbb{E}[q(t)E_t(\delta^t)|q(t)] + V\mathbb{E}[L_t(\delta^t)|q(t)] \\ &= B - Qq(t) + \mathbb{E}[q(t)E_t(\delta^t) + VL_t(\delta^t)|q(t)] \end{aligned} \quad (18)$$

where $L_t(\delta^t) \triangleq \sum_{k=1}^K [(1 - \delta_k^t) l_{k,rl}^{t,no} + \delta_k^t l_{k,rl}^{t,c}]$ is the response time at time slot t given the caching decision δ^t . $V(> 0)$ is a control parameter to adjust the tradeoff between energy consumption control and response time minimization. Now our attention has shifted from $P1$ to the minimization of the supremum of the *drift-plus-penalty* term, i.e., the right hand side of (18). Namely, the caching decision δ^t can be determined for time slot t by solving the following problem $P2$:

$$(P2) \quad \min_{\forall t, \delta^t} \{q(t)E_t(\delta^t) + VL_t(\delta^t)\} \quad (19)$$

$$s.t. (13), (14), (15) \quad (20)$$

It is worth mentioning that the constraint (12) in $P1$ has been incorporated into problem $P2$ itself by means of $\Delta(q(t))$.

Theorem 1: Based on the drift-plus-penalty term, the caching decisions obtained by solving $P2$ over time slots $t \in \{0, \dots, T-1\}$ are approximately optimal caching decisions w.r.t. $P1$.

Proof If there is a solution to $P1$, then there exists a caching decision δ^{t*} and l^* , while satisfying 1) $E_t(\delta^{t*}) \leq Q$, and 2) $l^* = \min_{\delta^t, \forall t} \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \sum_{k=1}^K L_t(\delta^t)$, based on [29].

Therefore, we have:

$$\begin{aligned} & \Delta(q(t)) + V\mathbb{E}[L_t(\delta^t)|q(t)] \\ & \leq B - Qq(t) + \mathbb{E}[q(t)E_t(\delta^t)|q(t)] + V\mathbb{E}[L_t(\delta^t)|q(t)] \\ & = B + \mathbb{E}[q(t)(E_t(\delta^t) - Q)|q(t)] + V\mathbb{E}[L_t(\delta^t)|q(t)] \\ & = B + q(t)\mathbb{E}[(E_t(\delta^t) - Q)|q(t)] + V\mathbb{E}[L_t(\delta^t)|q(t)] \\ & \stackrel{\dagger}{\leq} B + Vl^* \end{aligned}$$

Since $\Delta(q(t)) + V\mathbb{E}[L_t(\delta^t)|q(t)] \leq B + q(t)\mathbb{E}[(E_t(\delta^t) - Q)|q(t)] + V\mathbb{E}[L_t(\delta^t)|q(t)]$, for $\forall t \in \{0, \dots, T-1\}$ and valid caching decision δ^t , the inequality still holds when substituting δ^t by δ^{t*} , i.e., $B + q(t)\mathbb{E}[(E_t(\delta^{t*}) - Q)|q(t)] + V\mathbb{E}[L_t(\delta^{t*})|q(t)]$. Owing to $E_t(\delta^{t*}) - Q \leq 0$, the inequality (\dagger) holds.

Then we calculate the expectation of the above inequality and then perform a sum of the expectation over the time slots $t \in \{0, \dots, T-1\}$, namely,

$$\begin{aligned} & \sum_{t=0}^{T-1} \mathbb{E}[\Delta(q(t)) + V\mathbb{E}[L_t(\delta^t)|q(t)]] \\ & = \sum_{t=0}^{T-1} \mathbb{E}[\mathbb{E}[L(q(t+1)) - L(q(t))|q(t)] + V\mathbb{E}[L_t(\delta^t)|q(t)]] \\ & = \sum_{t=0}^{T-1} \mathbb{E}[L(q(t+1)) - L(q(t))] + V\mathbb{E}[L_t(\delta^t)] \\ & = \mathbb{E}[L(q(T)) - L(q(0))] + \sum_{t=0}^{T-1} V\mathbb{E}[L_t(\delta^t)] \\ & = \mathbb{E}[L(q(T))] + \sum_{t=0}^{T-1} V\mathbb{E}[L_t(\delta^t)] \leq \sum_{t=0}^{T-1} \mathbb{E}[B + Vl^*] \\ & = (B + Vl^*) * T \end{aligned}$$

Since $\mathbb{E}[L(q(T))] \geq 0$, we have

$$\sum_{t=0}^{T-1} V\mathbb{E}[L_t(\delta^t)] \leq (B + Vl^*) * T$$

Then,

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[L_t(\delta^t)] \leq l^* + \frac{B}{V}$$

Namely,

$$\frac{1}{T} \sum_{t=0}^{T-1} L_t(\delta^t) \leq l^* + \frac{B}{V}$$

This is because $L_t(\delta^t)$ is the expected response time, and according to the law of iterated expectations, $\mathbb{E}[L_t(\delta^t)] =$

$L_t(\delta^t)$. Thus, the caching decision obtained by solving $P2$ can be infinitely close to the solution of $P1$ by adjusting the variable V . \square

Remark: Let $\mathcal{G}(V)$ represent the optimality gap which denotes the difference between l^* and the solution of $P2$. The supremum of $P2$ is $l^* + B/V$, so $\mathcal{G}(V) = B/V$ and $O(\mathcal{G}(V)) = O(l^* + B/V) = O(1/V)$. As a consequence, if problem $P2$ can be solved within each time slot, $\mathcal{G}(V)$ can be bounded by $O(1/V)$.

Theorem 2: The caching decisions obtained by solving problem $P2$ over time slot $t \in \{0, \dots, T-1\}$ make inequality (12) always true.

Proof Let δ^{t*} denote the caching decision obtained by solving problem $P2$, i.e., the right-hand side of drift-plus-penalty inequality (18) in time slot t . Thus, we have:

$$\begin{aligned} & \Delta(q(t)) + V\mathbb{E}[L_t(\delta^{t*})|q(t)] \\ & \leq B - Qq(t) + \mathbb{E}[q(t)E_t(\delta^{t*})|q(t)] + V\mathbb{E}[L_t(\delta^{t*})|q(t)] \\ & \leq B - Qq(t) + VL_t^\dagger + q(t)(Q - \epsilon) \\ & = B - \epsilon q(t) + VL_t^\dagger \end{aligned}$$

where $\delta^{t\dagger}$ is any other valid caching decision in time slot t except δ^{t*} . L_t^\dagger is the corresponding response time with caching decision equal to $\delta^{t\dagger}$. Since $\mathbb{E}[E_t(\delta^{t*})] = E_t(\delta^{t*}) \leq Q$ [15], there exists a small enough $\epsilon > 0$, satisfying $E_t(\delta^{t*}) - Q \leq -\epsilon$, namely, $E_t(\delta^{t*}) \leq Q - \epsilon$.

By taking expectations of the above inequality, we have: $\mathbb{E}\{\Delta(q(t)) + V\mathbb{E}[L_t(\delta^{t*})|q(t)]\} \leq B - \epsilon\mathbb{E}[q(t)] + VL_t^\dagger$, namely,

$$\begin{aligned} & \mathbb{E}[L(q(t+1))] - \mathbb{E}[L(q(t))] + V\mathbb{E}[L_t(\delta^{t*})] \\ & \leq B - \epsilon\mathbb{E}[q(t)] + VL_t^\dagger \end{aligned}$$

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

$$\begin{aligned} & \mathbb{E}[L(q(T))] - \mathbb{E}[L(q(0))] + V \sum_{t=0}^{T-1} \mathbb{E}[L_t(\delta^{t*})] \\ & \leq BT + VL_t^\dagger T - \epsilon\mathbb{E}[q(t)]T \\ & \leq BT + VL_t^\dagger T = B'T \end{aligned}$$

where $B' = B + VL_t^\dagger$. Since $V \sum_{t=0}^{T-1} \mathbb{E}[L_t(\delta^{t*})] \geq 0$, thus $\mathbb{E}[L(q(T))] - \mathbb{E}[L(q(0))] \leq BT + VL_t^\dagger T$. Substitute $L(q(T))$ by $L(q(T)) = \frac{1}{2}q^2(T)$, we have:

$$\frac{1}{2}\mathbb{E}[q^2(T)] \leq B'T + \frac{1}{2}\mathbb{E}[q^2(0)] = B'T$$

According to the Cauchy inequality, i.e., $(\sum_{i=1}^n x_i y_i)^2 \leq (\sum_{i=1}^n x_i^2)(\sum_{i=1}^n y_i^2)$, we have,

$$(\mathbb{E}[q(T)])^2 \leq \mathbb{E}[q^2(T)] \leq 2B'T$$

Thus, $\mathbb{E}[q(T)] \leq \sqrt{2B'T}$, so

$$\frac{1}{T}\mathbb{E}[q(T)] \leq \sqrt{\frac{2B'}{T}}$$

Base on Lemma 1, we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[E_t(\delta^t) - Q] \leq \frac{1}{T} \mathbb{E}[q(T)] \leq \sqrt{\frac{2B'}{T}}$$

Therefore,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[E_t(\delta^t) - Q] \leq \lim_{T \rightarrow \infty} \sqrt{\frac{2B'}{T}} = 0$$

Namely,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} [E_t(\delta^t) - Q] \leq 0$$

As a consequence, we can see that 1) the queue backlog is mean rate stable; and 2) the caching decision obtained by solving $P2$ can solve $P1$ without constraint violation. \square

Based on the above descriptions, we have proposed a Lyapunov based online caching decision (LOCD) algorithm to obtain approximately optimum solution to problem $P1$, as shown in Alg. 1. Specifically, LOCD seeks the optimal solution along each time slot. Within each time slot, given per-slot constraints on the expected response time and energy consumption (i.e., L_t^{max} and E_t^{max}), respectively, we strive to achieve approximately optimum solution by solving problem $P2$ (lines 3–5). Then, the per-slot best caching decision can be obtained, i.e., δ^{t*} . LOCD calculates the optimal response time based on δ^{t*} and then the per-slot response time is accumulated (lines 7–8). Finally, the mean value of the objective function (i.e., $avg = L/T$) can be retrieved.

Algorithm 1: Lyapunov based Online Caching Decision (LOCD)

Input: $q(0), Q, \mathcal{C}, \mathcal{A}, V, T$
Output: Optimum solution to $P1$

```

1  $L = 0$ ;
2 for  $t = 0$  to  $T - 1$  do
3   Observe and record  $\lambda^t, r^t$ ;
4   Set  $L_t^{max}$  and  $E_t^{max}$ ;
5   Obtain  $\delta^{t*}$  by solving problem  $P2$ ;
6    $L = L + L_t(\delta^{t*})$ ;
7   Calculate  $E_t(\delta^t)$  based on Eq. 10;
8    $q(t+1) = \max[q(t) - Q, 0] + E_t(\delta^t)$ ;
9 end
10  $avg = L/T$ ;
11 return  $avg$ ;
```

B. Algorithm Design for Solving $P2$

Let $\delta^{t*} = \arg \min_{\delta^t, \forall t} \{q(t)E_t(\delta^t) + VL_t(\delta^t)\}$ while satisfying the constraints (13), (14), (15). Then δ^{t*} can be viewed as the best caching decision in time slot t in terms of $L_t(\delta^t)$. Given $q(t)$ and V , the problem $P2$ is actually a 0-1 knapsack problem. Owing to its well-known NP-hardness, exhaustive search takes exponential time and thus makes itself infeasible in reality. Evolutionary algorithms such as GA and PSO [14], [20] can be adopted to solve it. However, these iteration-based

methods are usually time-consuming and thus impracticable in time-sensitive scenarios. Accordingly, a heuristic algorithm suitable for time-sensitive scenario is required to solve $P2$. Let $\mathcal{F} = q(t)E_t(\delta^t) + VL_t(\delta^t)$ and by substituting $E_t(\delta^t)$ and $L_t(\delta^t)$ respectively, we have:

$$\begin{aligned} \mathcal{F}(\delta^t) &= q(t)E_t(\delta^t) + VL_t(\delta^t) \\ &= \sum_{k=1}^K [q(t)e_k^{t,no} + Vl_{k,rl}^{t,no}] - \sum_{k=1}^K \delta_k^t [q(t)\gamma_k + V(l_{k,rl}^{t,no} - l_{k,rl}^{t,c})] \\ &\triangleq \mathcal{P} - \mathcal{G} \end{aligned} \quad (21)$$

where $\mathcal{P} = \sum_{k=1}^K [q(t)e_k^{t,no} + Vl_{k,rl}^{t,no}]$, independent of the caching decision δ^t , can be determined within each time slot, and $\mathcal{G} = \sum_{k=1}^K \delta_k^t [q(t)\gamma_k + V(l_{k,rl}^{t,no} - l_{k,rl}^{t,c})]$ depends upon δ^t . To minimize \mathcal{F} means to maximize \mathcal{G} in essence, since \mathcal{P} is unchangeable given the time slot t .

To meet the strict response time requirement, a greedy algorithm is proposed to minimize \mathcal{F} . In particular, a greedy heuristic is given in what follows. Let $g_k^t = q(t)\gamma_k + V(l_{k,rl}^{t,no} - l_{k,rl}^{t,c})$, and the applications in \mathcal{A} are sorted in the descending order of g_k^t/d_k . The corresponding greedy algorithm is shown in Alg. 2. A list \mathcal{L} is used to store the sorted applications. S_0 and S_1 store the indexes of tasks that are cached and not cached, respectively. Therefore, the caching decision profile δ^t can be easily constructed based on S_0 and S_1 . Initially, $S_1 = \emptyset$ and $S_0 = \{0, 1, \dots, K-1\}$, respectively. We check application a_k in \mathcal{L} one by one. If the sum of input data of all already cached applications plus current a_k does not exceed the constraint \mathcal{C} , we tend to cache a_k . However, the per-slot energy and latency constraints are examined respectively before updating the decision caching profile (lines 12–15). In the meanwhile, the accumulated energy consumption and response time are obtained, respectively (lines 10–11). If the sum of input data of all already cached applications plus current a_k exceeds the constraint \mathcal{C} , the algorithm checks the application after a_k in \mathcal{L} . In the meanwhile, the accumulated energy consumption and response time are also calculated. Based on Eq. 21, the minimal \mathcal{F} can be retrieved. Finally, S_0 and S_1 are returned.

Remark: The time of GASP is mainly taken to sort the applications in \mathcal{A} . Various sort algorithms can be utilized to achieve this goal with the time complexity of $O(K \log K)$, where K is the number of applications. In addition, the time taken to check applications in \mathcal{A} is $O(K)$. Thus, the time complexity of GASP is $O(K \log K) + O(K) = O(K \log K)$, which is much better than the evolutionary algorithms. GASP endeavors to obtain the approximate optimal solution according to the greedy heuristic. However, similar to general 0-1 knapsack problem solved by greedy algorithms, the drawback of GASP is that the vacant part of knapsack can depreciate the value of knapsack if the knapsack is not fully occupied. Therefore, GASP needs to be enhanced for the performance improvement.

C. Enhanced Algorithm for Solving Problem $P2$

To narrow down the optimality gap between the approximate optimal value and true optimal value, in this section

Algorithm 2: Greedy Algorithm for Solving P2 (GASP)

Input: $q(t)$, $e_k^{t,no}$, $l_{k,rl}^{t,c}$, $l_{k,rl}^{t,no}$, C , \mathcal{A} , V , γ_k
Output: S_0 and S_1

```

1  $S_0 = \{0, 1, \dots, K-1\}$ ;
2  $S_1 = \emptyset$ ;
3  $g_k^t = q(t)\gamma_k + V(l_{k,rl}^{t,no} - l_{k,rl}^{t,c})$ , for  $\forall k \in [1, K]$ ;
4  $\mathcal{L}$  stores the applications in descending order of  $g_k^t/d_k$ ;
5  $s = 0, L = 0, E = 0$ ;
6 for each  $a_k$  in  $\mathcal{L}$  do
7    $s = s + d_k$ ;
8   if  $s \leq C$  then
9     Calculate  $l_{k,rl}^{t,c}$  and  $e_k^{t,c}$ ;
10     $L = L + l_{k,rl}^{t,c}$ ;
11     $E = E + e_k^{t,c}$ ;
12    if  $L \leq L_t^{max}$  and  $E \leq E_t^{max}$  then
13       $\delta_k^t = 1$ ;
14       $S_1 = S_1 \cup \{k\}$ ;
15       $S_0 = S_0 \setminus \{k\}$ ;
16    else
17       $L = L - l_{k,rl}^{t,c}$ ;
18       $E = E - e_k^{t,c}$ ;
19    end
20  else
21     $s = s - d_k$ ;
22    Calculate  $l_{k,rl}^{t,no}$  and  $e_k^{t,no}$ ;
23     $L = L + l_{k,rl}^{t,no}$ ;
24     $E = E + e_k^{t,no}$ ;
25  end
26 end
27  $\mathcal{F}_{min} = q(t)E + VL$ ;
28 return  $S_0$  and  $S_1$ ;

```

we apply two other greedy rules to the approximate solution obtained by GASP.

Definition 1: One-dimensional increment $\Delta_k(\delta^t)$ is defined as $\Delta_k(\delta^t) = \mathcal{G}(\delta_1^t, \dots, 1_{(k)}, \dots, \delta_K^t) - \mathcal{G}(\delta_1^t, \dots, 0_{(k)}, \dots, \delta_K^t) = q(t)(e_k^{t,no} - e_k^{t,c}) + V(l_{k,rl}^{t,no} - l_{k,rl}^{t,c}) = q(t)\gamma_k + V(l_{k,rl}^{t,no} - l_{k,rl}^{t,c})$.

$\Delta_k(\delta^t)$ can be used to evaluate the optimum increment to the objective function \mathcal{G} by unilaterally changing δ_k^t from 0 to 1 while keeping other caching decisions $\{\delta_i^t | 1 \leq i \leq K, i \neq k\}$ unchanged.

Definition 2: Two-dimensional increment $\Delta_{kl}(\delta^t)$ is defined as $\Delta_{kl}(\delta^t) = \mathcal{G}(\delta_1^t, \dots, 1_{(k)}, \dots, 0_{(l)}, \dots, \delta_K^t) - \mathcal{G}(\delta_1^t, \dots, 0_{(k)}, \dots, 1_{(l)}, \dots, \delta_K^t) = q(t)(\gamma_k - \gamma_l) + V(l_{k,rl}^{t,no} - l_{k,rl}^{t,c} - l_{l,rl}^{t,no} + l_{l,rl}^{t,c})$.

$\Delta_{kl}(\delta^t)$ can denote the optimum increment to \mathcal{G} by swapping the caching decisions of two services (e.g., service k from S_1 and l from S_0 , respectively), with one varying from 0 to 1 and the other from 1 to 0. If $\Delta_{kl}(\delta^t) < 0$, the objective function \mathcal{F} can be further optimized by exchanging the caching decisions of two tasks from S_1 and S_0 , respectively.

Accordingly, we can improve the performance of GASP with the aid of $\Delta_k(\delta^t)$ and $\Delta_{kl}(\delta^t)$. To be specific, given a decision profile δ^t obtained by GASP, i.e., S_0 and S_1 , the procedure of enhancing GASP, denoted by EnGASP, is shown

in Alg. 3. EnGASP mainly consists of two steps, with one called the filling stage (lines 1–11) and the other called the swapping stage (lines 13–21). The filling stage, targeted at the vacant part of knapsack, strives to fill the knapsack with uncached applications from S_0 . The swapping stage on the other hand aims to find a better caching decision profile by swapping two applications from S_1 and S_0 respectively. For example, $\Delta_{k^*l^*}(\delta^t)$ smaller than 0 means that \mathcal{G} can be increased by exchanging the caching decisions of applications a_{k^*} and a_{l^*} . In the meanwhile, if the caching size constraint can be satisfied after the exchanging, \mathcal{F} can be further optimized by increasing \mathcal{G} .

Algorithm 3: Enhanced Greedy Algorithm for Solving P2 (EnGASP)

Input: \mathcal{F}_{min} , δ^t , S_0 , S_1
Output: δ^t

```

1 while  $S_0 \neq \emptyset$  do
2   Obtain  $k^*$  by solving:
3    $k^* = \arg \max \{\Delta_k(\delta^t) | k \in S_0\}$ ;
4   if  $\sum_i \{d_i | i \in S_1\} + d_{k^*} \leq C$  then
5     Update  $\mathcal{F}_{min}$ ;
6      $\delta_{k^*}^t = 1$ ;
7      $S_1 = S_1 \cup \{k^*\}$ ;
8      $S_0 = S_0 \setminus \{k^*\}$ ;
9   else
10     $S_0 = S_0 \setminus \{k^*\}$ ;
11  end
12  $S_0 = \{1, \dots, K\} \setminus S_1$ ;
13 while  $S_0 \neq \emptyset$  and  $S_1 \neq \emptyset$  do
14   Obtain  $(k^*, l^*)$  by solving:
15    $(k^*, l^*) = \arg \min \{\Delta_{kl}(\delta^t) | k \in S_1, l \in S_0\}$ ;
16   if  $\Delta_{k^*l^*}(\delta^t) < 0$  and
17    $\sum_i \{d_i | i \in S_1 \setminus \{k^*\} \cup \{l^*\}\} \leq C$  then
18     Update  $\mathcal{F}_{min}$ ;
19      $\delta_{k^*}^t = 0$  and  $\delta_{l^*}^t = 1$ ;
20      $S_1 = S_1 \setminus \{k^*\}$ ;
21      $S_0 = S_0 \setminus \{l^*\}$ ;
22  end
23 return  $\delta^t$ ;

```

Remark: It is worth mentioning that EnGASP attempts to enhance GASP rather than replace it, since the input parameters of EnGASP are the output parameters of GASP. Intuitively, the performance of EnGASP is at least as good as GASP. However, as for the time complexity of EnGASP, we can see that the main two loops (i.e., line 1 and line 13) take $O(K)$ and $O(K)$, respectively. To obtain k^* usually needs to traverse the set S_0 which takes the time of $O(K)$. On the other hand, to obtain (k^*, l^*) usually requires to find the pair of (l, k) by traversing S_0 and S_1 , which takes time of $O(|S_1|)$ and $O(|S_0|)$, respectively. Thus, the time to find (k^*, l^*) is $O(|S_0| \cdot |S_1|)$. Since $O(|S_0| \cdot |S_1|) \leq O((|S_0| + |S_1|)/2)^2 = O(K^2/4) = O(K^2)$, the time complexity of EnGASP can be derived as $O(K) \cdot O(K) + O(K) \cdot O(K^2) = O(K^3)$. Compared

TABLE II
PARAMETER SETTINGS

Parameter	Value	Parameter	Value
T	[0, 1000]	K	20
d_k	[1, 50]	s_k	[40, 70]
C	200	r_k^t	150
β_k^t	[1, 3]	λ_k^t	[100, 200]
π	400	f_e^t	[1000, 2000]
m	3	$\kappa\varepsilon$	1e-5
Q	1.2	E_t^{max}	20
L_t^{max}	100	V	[300, 1000]
γ_k^t	(0, 0.1)	$d_{k,init}^{t,no}$	(0, 1)

to the iteration-based evolutionary algorithms [14], [20], both GASP and EnGASP are of polynomial time. Therefore, both are really appropriate for time-sensitive scenarios in VEC systems. In addition, extensive simulation is still needed to evaluate EnGASP and GASP w.r.t. efficiency and effectiveness.

V. NUMERIC EVALUATION

Extensive experiments have been carried out to evaluate the proposed caching strategy in terms of response time and energy consumption in this section. In what follows, we will report the experimental settings and simulation results, respectively.

A. Experimental Settings

The involved parameters with the corresponding values are shown in Table II. Although the involved parameters in the simulation are given manually, our previous works [14], [20] can bring rich experience to the parameter settings, including the duration for each time slot π , $\kappa\varepsilon$, f_e^t and so on. For instance, π should be set appropriately such that on one hand the edge server S has enough time to make caching decisions, and on the other hand tasks offloaded to S can be completely accomplished by the end of time slot. The number of time slots varies from 0 to 1000 and the number of applications is set to 20. Within each time slot, these vehicular applications are randomly generated according to d_k and s_k . For the caching related parameters, the global energy constraint Q across different time slots is set to 1.2. For simplicity, we assume that the maximal per-slot energy consumption constraint is the same for each time slot (i.e., 20) in the simulation. The control parameter for the drift-plus-penalty term V ranges from 300 to 1000.

B. Simulation Results

First, we evaluate the performance of GASP and EnGASP with different time slots respectively. The experimental result is shown in Fig. 3, where the x-coordinate denotes the time slots and the y-coordinate denotes the objective values corresponding to the problem $P2$. From this figure, several observations can be revealed as follows. First, as theoretically analyzed earlier, the performance of EnGASP is at least as good as GASP. The experimental result conforms to the

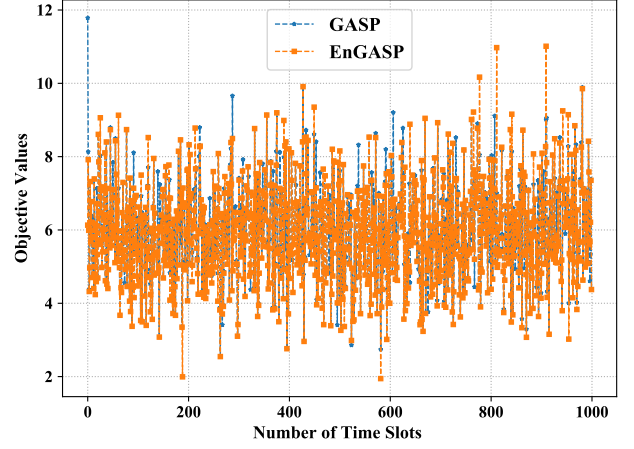


Fig. 3. Performance comparison with different time slots

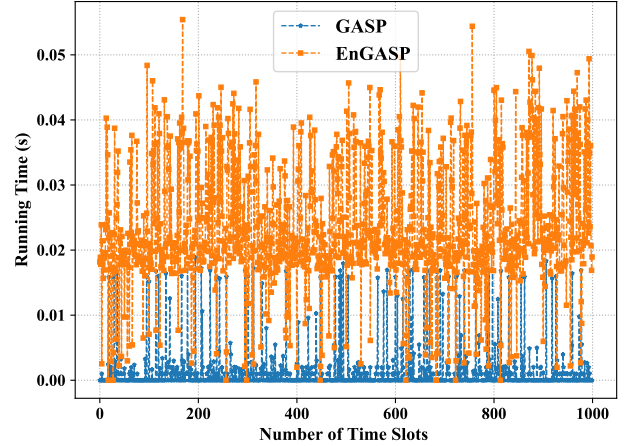


Fig. 4. Response time comparison with different time slots

theoretical analysis. Second, the performance of GASP at time slot 0 is the worst compared to other cases, which however is comprehensible and explicable as follows. According to Eq. 16, the calculation of $q(t+1)$ at time slot t depends on $q(t)$ and $E_t(\delta^t)$. At the beginning of time slot (i.e., $t = 0$), we assume that all the vehicular applications are not cached at S . Thus, the objective value can be very large. Third, EnGASP helps improve the performance of GASP by two main stages (i.e., the filling step and the swapping step). Across different time slots, the two steps further optimize GASP to a certain extent.

Fig. 4 shows the change of the running time of GASP and EnGASP as the number of time slots increases. First, it is very clear that the running time of EnGASP is much larger than that of GASP. Most of the time, GASP can make a caching decision in real time, while EnGASP can make it within tens of milliseconds. It is understandable since EnGASP needs the output of GASP as input parameters, i.e., EnGASP runs based on the caching decision obtained from GASP. As a result, the response time of EnGASP includes that of GASP, and thus it is of polynomial time. Second, both GASP and EnGASP need

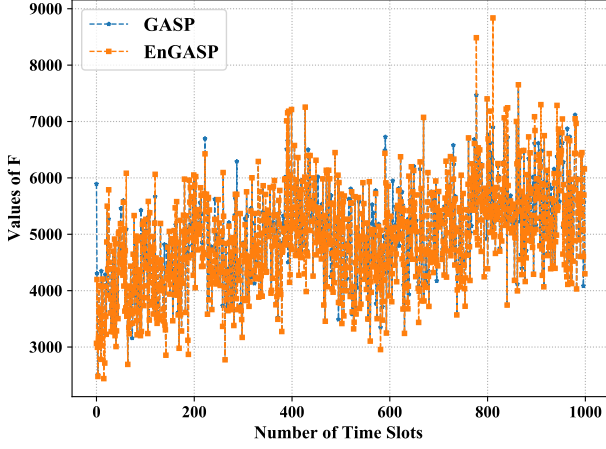


Fig. 5. \mathcal{F} minimization with two heuristic rules

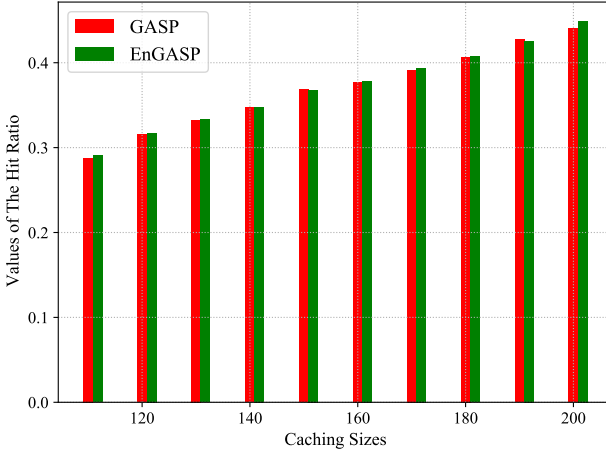


Fig. 6. Hit ratio comparison with different caching sizes

complete accomplishment within each time slot. Furthermore, the tasks are generated randomly within each time slot which means they are independent at different time slots. Therefore, the running time of both GASP and EnGASP do not increase as the number of time slots increases. Third, both GASP and EnGASP can make a caching decision in almost real time. It is of important significance to satisfy the strict time requirement of the vehicular applications.

Long-term energy constraint Q is converted into the per-slot energy constraint by Lyapunov optimization technology, so we can pay attention to the per-slot problem optimization. GASP and EnGASP try to obtain the best caching decisions w.r.t. problem $P1$ by minimizing \mathcal{F} . The experiment has been conducted to evaluate the performance when minimizing \mathcal{F} and the simulation result is shown in Fig. 5, where the x-coordinate represents the number of time slots and y-coordinate represents the values of \mathcal{F} . It can be easily observed that EnGASP is better than GASP when the number of time slots is the same. Generally speaking, no matter how the number of time slots varies, EnGASP is better than GASP in most cases. On the other hand, the situation similar to the case shown in Fig. 3

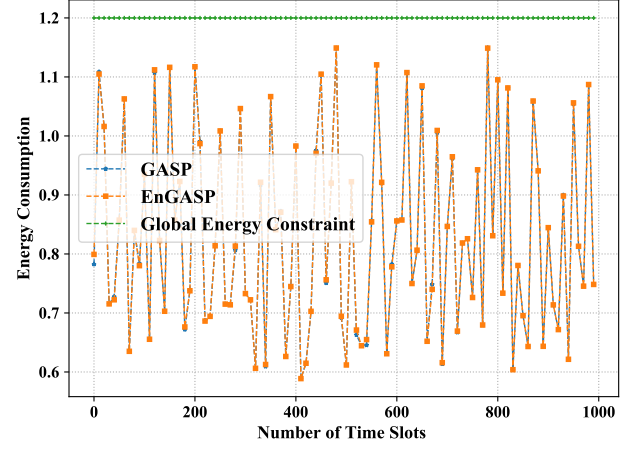


Fig. 7. The average energy consumption compared to the global energy constraint

happens when the time slot $t = 0$ and the reason is also the same as that in Fig. 3.

The metric hit ratio is an important performance indicator to evaluate the performance of the caching strategies in the information-centric networks. Furthermore, it is still applicable for evaluating the application-oriented caching strategies. In particular, Fig. 6 shows the variation of the hit ratio when GASP and EnGASP are applied under different caching sizes. In this experiment, the caching sizes range from 110 to 200 with a step of 10. It is obviously that the hit ratio increases for both GASP and EnGASP, as the caching sizes increase, this is because larger caching capabilities will enable more services to be cached. As expected, EnGASP is slightly better than GASP w.r.t. the hit ratio in most cases.

Another set of experiments is conducted to evaluate the average energy consumption at S across different time slots. More importantly, we need to check whether there are any constraint violations of energy consumption at different time slots. The simulation results are shown in Fig. 7. First of all, no matter how the number of time slots increases, the energy consumptions at S vary roughly from 0.5 to 1.2. The global energy constraint Q is 1.2, so there is no energy violation at all for the arbitrary number of time slots. Therefore, the experimental results accord with the theoretical analysis. Second, the energy consumption at S within each time slot depends upon the current caching decision profile. Generally, the more the number of applications which are cached at S , the more the energy consumptions at S , for the reason that given a certain period of time, the running state of virtual environments should be maintained so as to avoid additional time overheads on the virtual initialization. The difference in the caching decision profile at different time slots leads to the different energy consumptions at different time slots.

As discussed earlier, some iteration-based evolutionary algorithms can also be adopted to obtain the optimal values of \mathcal{F} . In particular, we investigate the performance of these algorithms (e.g., GA and PSO) w.r.t., the capabilities to obtain the optimal values of \mathcal{F} and the time complexity. The simulation results are shown in Fig. 8, 9 10, and 11, respectively.

Fig. 8 and Fig. 9 show the optimal values and the running time of the four approaches, respectively, when the number of time slots increases. Obviously, at least two conclusions can be drawn based on the observations. First, GA and PSO indeed have better capabilities to obtain the optimal values of \mathcal{F} compared to our approaches GASP and EnGASP. Since the offloading requests are generated randomly in each time slot, and thus they are independent in different time slots. The optimization gap between the evolutionary algorithms and our approaches is not becoming wider as the number of time slots increases. Second, the running times of GA and PSO are much longer than GASP and EnGASP. For example, it averagely takes about 2 seconds for GA to obtain the optimal value of \mathcal{F} , while it averagely tasks about 1.65 seconds for PSO to obtain the value of \mathcal{F} . GASP and EnGASP on the other hand can achieve the realtime response no matter how the number of time slots changes.

When we compare the four approaches under different number of applications, GASP and PSO are hardly acceptable despite the advantages over our approaches. As denoted in Fig. 11, both of GA and PSO take seconds to obtain the better values of \mathcal{F} . As a result, we can carefully draw a conclusion that the evolutionary algorithms are not as good as our approaches in the time-sensitive application scenario emphasized by this paper. Our heuristic rules cater for the strict time requirements of vehicular applications at the expense of the precision of the solution, which however is acceptable to a great extent. In contrast, GA and PSO may display better performance with regards to precision, but they usually take sub-seconds or even seconds to achieve this goal. Such amount of time is unacceptable for the time-sensitive vehicular applications.

Next, we investigate the effect of V on the objective values obtained in LOCD. To be specific, GASP and EnGASP are adopted to solve problem P_2 (line 5 in LOCD), respectively. As shown in Theorem 1, the caching decision obtained by solving P_2 can be infinitely close to the solution to P_1 by adjusting the variable V . The value of V varies from 800 to 1000, and the simulation results are shown in Fig. 12. From the figure, we can observe that LOCD with EnGASP has significant advantages over LOCD with GASP. On the other hand, owing to the random generation of tasks at different time slots, the best objective values are independent of each other across different time slots. Accordingly, it makes sense that the best objective values fluctuate a lot across different time slots.

The last set of experiments has been carried out to evaluate the performance of LOCD with the increasing number of time slots. Due to the global energy constraint Q , LOCD reveals its convergence either with GASP or EnGASP no matter how the number of time slots increases. For example, when the number of time slots is equal to 1000, the corresponding objective value is about 6.03. In spite of the fluctuation of these values across different time slots, this fluctuation has been confined to a limited interval (e.g., between 5.7 and 6.35). Furthermore, LOCD with EnGASP still outstands that with GASP all the time.

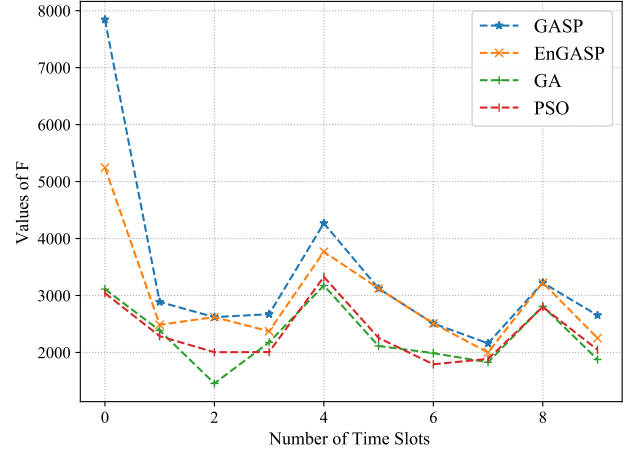


Fig. 8. Performance comparison with different number of time slots

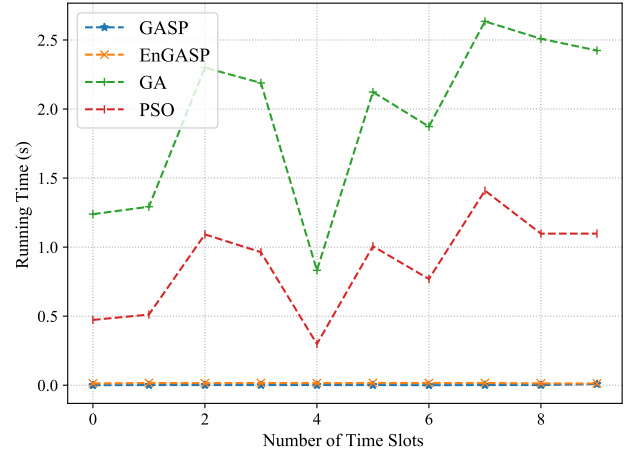


Fig. 9. Running time comparison with different number of time slots

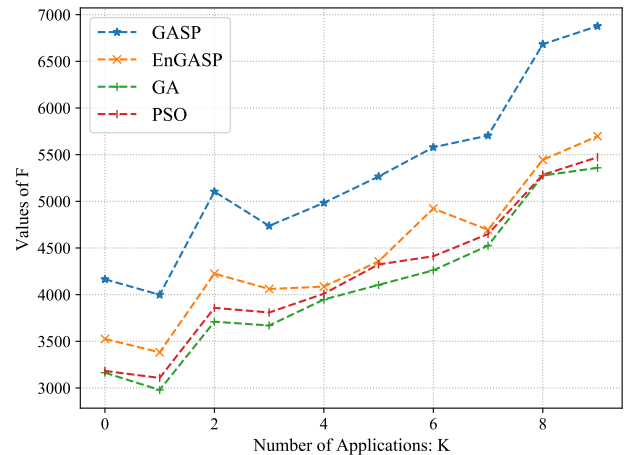


Fig. 10. Performance comparison with different number of applications

VI. CONCLUSION

VEC brings considerable benefits for vehicular applications in smart transportation. However, the performance of VEC is

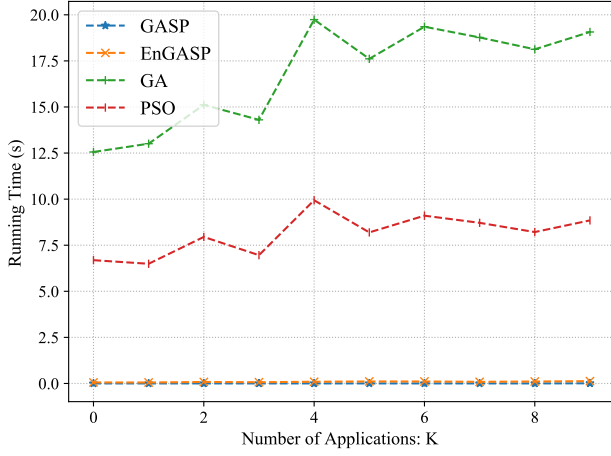


Fig. 11. Running time comparison with different number of applications

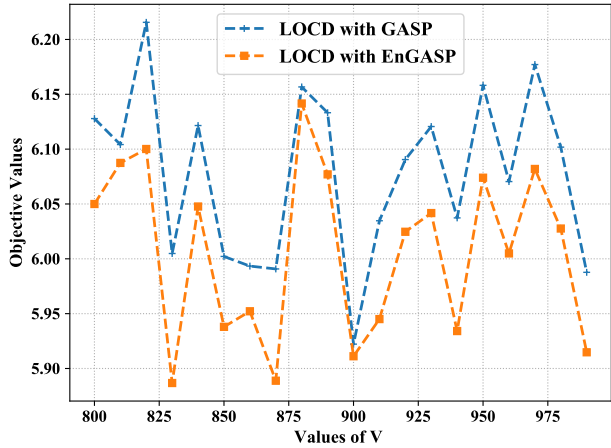


Fig. 12. The effect of V on the objective values

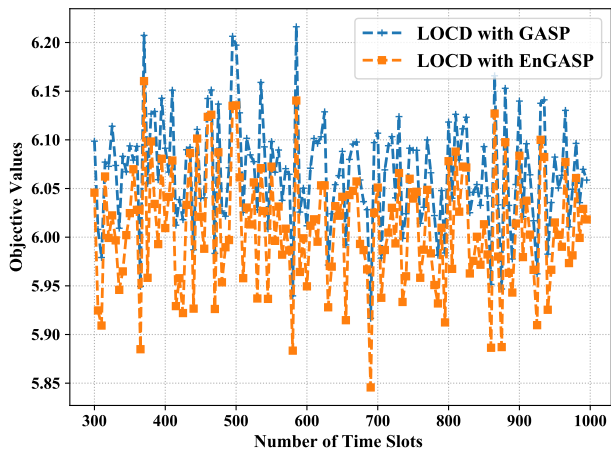


Fig. 13. Convergence investigation with different number of time slots

still challenged by high mobility of vehicles and limited communication resources, especially considering that the response time acts as the optimization objective during application outsourcing. In this paper, we apply application caching to VEC

to optimize the response time for the outsourced applications while satisfying the time slot spanned energy consumption. The Lyapunov optimization technology is adopted to convert the global energy constraint into per-slot energy constraints, so as to facilitate the response time optimization. Furthermore, two greedy heuristics are incorporated into the drift-plus-penalty based algorithm for helping find the approximate optimal solution. We have evaluated the approach via a series of experiments. The simulation results reveal the advantages of our algorithms in terms of response time and energy consumption.

REFERENCES

- [1] J. Wang, C. Jiang, K. Zhang, T. Q. S. Quek, Y. Ren, and L. Hanzo, "Vehicular sensing networks in a smart city: Principles, technologies and applications," *IEEE Wirel. Commun.*, vol. 25, no. 1, pp. 122–132, 2018.
- [2] N. Zhang, S. Zhang, P. Yang, O. Alhussein, W. Zhuang, and X. S. Shen, "Software defined space-air-ground integrated vehicular networks: Challenges and solutions," *IEEE Commun. Mag.*, vol. 55, no. 7, pp. 101–109, 2017.
- [3] J. Zhang, H. Guo, J. Liu, and Y. Zhang, "Task offloading in vehicular edge computing networks: A load-balancing solution," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 2, pp. 2092–2104, 2020.
- [4] Y. Yao, B. Xiao, W. Wang, G. Yang, X. Zhou, and Z. Peng, "Real-time cache-aided route planning based on mobile edge computing," *IEEE Wireless Communications*, vol. 27, no. 5, pp. 155–161, 2020.
- [5] G. Zhang, Y. Li, and T. Lin, "Caching in information centric networking: A survey," *Computer Networks*, vol. 57, no. 16, pp. 3128–3141, 2013.
- [6] M. Zhang, H. Luo, and H. Zhang, "A survey of caching mechanisms in information-centric networking," *IEEE Communications Surveys and Tutorials*, vol. 17, no. 3, pp. 1473–1499, 2015.
- [7] J. Xu, L. Chen, and P. Zhou, "Joint service caching and task offloading for mobile edge computing in dense networks," in *IEEE INFOCOM 2018 - IEEE Conference on Computer Communications*, 2018, pp. 207–215.
- [8] Y. Hao, M. Chen, L. Hu, M. S. Hossain, and A. Ghoniem, "Energy efficient task caching and offloading for mobile edge computing," *IEEE Access*, pp. 11 365–11 373, 2018.
- [9] H. Xing, J. Cui, Y. Deng, and A. Nallanathan, "Energy-efficient proactive caching for fog computing with correlated task arrivals," in *2019 IEEE 20th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 2019, pp. 1–5.
- [10] L. Yang, J. Cao, G. Liang, and X. Han, "Cost aware service placement and load dispatching in mobile cloud systems," *IEEE Transactions on Computers*, vol. 65, no. 5, pp. 1440–1452, 2016.
- [11] G. Avino, M. Malinverno, F. Malandrino, C. Casetti, and C. F. Chiasserini, "Characterizing docker overhead in mobile edge computing scenarios," in *Proceedings of the Workshop on Hot Topics in Container Networking and Networked Systems*, ser. HotConNet '17. Association for Computing Machinery, 2017, pp. 30–35.
- [12] X. Hou, Z. Ren, J. Wang, W. Cheng, Y. Ren, K. Chen, and H. Zhang, "Reliable computation offloading for edge-computing-enabled software-defined iov," *IEEE Internet Things J.*, vol. 7, no. 8, pp. 7097–7111, 2020.
- [13] S. Li, S. Lin, L. Cai, W. Li, and G. Zhu, "Joint resource allocation and computation offloading with time-varying fading channel in vehicular edge computing," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 3, pp. 3384–3398, 2020.
- [14] C. Tang, C. Zhu, X. Wei, H. Wu, Q. Li, and J. J. P. C. Rodrigues, "Intelligent resource allocation for utility optimization in rsu-empowered vehicular network," *IEEE Access*, vol. 8, pp. 94 453–94 462, 2020.
- [15] H. Wu, Y. Sun, and K. Wolter, "Energy-efficient decision making for mobile cloud offloading," *IEEE Transactions on Cloud Computing*, vol. 8, no. 2, pp. 570–584, 2020.
- [16] Z. Qin, S. Leng, J. Zhou, and S. Mao, "Collaborative edge computing and caching in vehicular networks," in *2020 IEEE Wireless Communications and Networking Conference (WCNC)*, 2020, pp. 1–6.
- [17] Y. Dai, D. Xu, K. Zhang, S. Maharjan, and Y. Zhang, "Deep reinforcement learning and permissioned blockchain for content caching in vehicular edge computing and networks," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 4, pp. 4312–4324, 2020.

- [18] X. Hou, Y. Li, M. Chen, D. Wu, D. Jin, and S. Chen, "Vehicular fog computing: A viewpoint of vehicles as the infrastructures," *IEEE Transactions on Vehicular Technology*, pp. 3860–3873, 2016.
- [19] M. Laroui, B. Nour, H. Mounsla, H. Afifi, and M. A. Cherif, "Mobile vehicular edge computing architecture using rideshare taxis as a mobile edge server," in *2020 IEEE 17th Annual Consumer Communications Networking Conference (CCNC)*, 2020, pp. 1–2.
- [20] C. Tang, X. Wei, C. Zhu, Y. Wang, and W. Jia, "Mobile vehicles as fog nodes for latency optimization in smart cities," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 9, pp. 9364–9375, 2020.
- [21] Y. Dai, D. Xu, S. Maharjan, and Y. Zhang, "Joint load balancing and offloading in vehicular edge computing and networks," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 4377–4387, 2019.
- [22] Z. Ning, P. Dong, X. Wang, J. J. P. C. Rodrigues, and F. Xia, "Deep reinforcement learning for vehicular edge computing: An intelligent offloading system," *ACM Trans. Intell. Syst. Technol.*, vol. 10, no. 6, pp. 60:1–60:24, 2019.
- [23] D. Chen, Y. Liu, B. Kim, J. Xie, C. S. Hong, and Z. Han, "Edge computing resources reservation in vehicular networks: A meta-learning approach," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 5, pp. 5634–5646, 2020.
- [24] N. Cheng, H. Zhou, L. Lei, N. Zhang, Y. Zhou, X. Shen, and F. Bai, "Performance analysis of vehicular device-to-device underlay communication," *IEEE Trans. Veh. Technol.*, vol. 66, no. 6, pp. 5409–5421, 2017.
- [25] J. Yao and N. Ansari, "Fog resource provisioning in reliability-aware iot networks," *IEEE Internet of Things Journal*, vol. 6, no. 5, pp. 8262–8269, 2019.
- [26] X. Chen, L. Jiao, W. Li, and X. Fu, "Efficient multi-user computation offloading for mobile-edge cloud computing," *IEEE/ACM Transactions on Networking*, vol. 24, no. 5, pp. 2795–2808, 2016.
- [27] X. Wei, C. Tang, J. Fan, and S. Subramaniam, "Joint optimization of energy consumption and delay in cloud-to-thing continuum," *IEEE Internet of Things Journal*, vol. 6, no. 2, pp. 2325–2337, 2019.
- [28] L. Georgiadis, M. J. Neely, and L. Tassiulas, "Resource allocation and cross-layer control in wireless networks," *Foundations & Trends in Networking*, vol. 1, no. 1, pp. 1–144, 2006.
- [29] 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.



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, 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, big data.



Chunsheng Zhu is an Associate Professor in the SUSTech Institute of Future Networks at Southern University of Science and Technology in China. He is also an Associate Researcher in the PCL Research Center of Networks and Communications at Peng Cheng Laboratory in China. He received the Ph.D. Degree in Electrical and Computer Engineering from The University of British Columbia, Canada. He has authored more than 100 publications published by refereed international journals (e.g., *IEEE Transactions on Industrial Electronics*, *IEEE Transactions*

on Computers, *IEEE Transactions on Information Forensics and Security*, *IEEE Transactions on Industrial Informatics*, *IEEE Transactions on Vehicular Technology*, *IEEE Transactions on Emerging Topics in Computing*, *IEEE Transactions on Cloud Computing*, *ACM Transactions on Embedded Computing Systems*, *ACM Transactions on Cyber-Physical Systems*), magazines (e.g., *IEEE Communications Magazine*, *IEEE Wireless Communications Magazine*, *IEEE Network Magazine*), and conferences (e.g., *IEEE INFOCOM*, *IEEE IECON*, *IEEE SECON*, *IEEE DCOSS*, *IEEE ICC*, *IEEE GLOBECOM*). His research interests mainly include Internet of Things, wireless sensor networks, cloud computing, big data, social networks, and security.



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 in the Center for Applied Mathematics, Tianjin University. His research interests include model-based evaluation, wireless and mobile network systems, mobile cloud computing and deep learning.



Qing Li is currently a Chair Professor (Data Science) and the Head of the Department of Computing, the Hong Kong Polytechnic University. Formerly, he was the founding Director of the Multimedia software Engineering Research Centre (MERC), and a Professor at City University of Hong Kong where he worked in the Department of Computer Science from 1998 to 2018. Prior to these, he has also at the Hong Kong University of Science and Technology and the Australian National University (Canberra, Australia). Prof. Li is currently a Fellow of IET/IEEE, a Senior Member of IEEE, a member of ACM-SIGMOD and IEEE Technical Committee on Data Engineering. He is the chairperson of the Hong Kong Web Society, and also served/is serving as an executive committee (EXCO) member of IEEE-Hong Kong Computer Chapter and ACM Hong Kong Chapter. In addition, he serves as a councilor of the Database Society of Chinese Computer Federation (CCF), a member of the Big Data Expert Committee of CCF, and is a Steering Committee member of DASFAA, ER, ICWL, UMEDIA, and WISE Society.



Joel J. P. C. Rodrigues (S'01-M'06-SM'06-F'20) is a professor at the Federal University of Piau, Brazil and a senior researcher at Instituto de Telecomunicações, Portugal. He is the leader of the Internet of Things Research Group (CNPq), Director for Conference Development - IEEE ComSoc Board of Governors, an IEEE Distinguished Lecturer, Technical Activities Committee Chair of the IEEE ComSoc Latin America Region Board, and Past Chair of the IEEE ComSoc eHealth and Communications Software TCs. He is the Editor-in-Chief of the International Journal of E-Health and Medical Communications and has authored and coauthored over 780 publications in refereed international journals and conferences, 3 books, 2 patents, and 1 ITU-T Recommendation.