

QoS-Aware Online Service Provisioning and Updating in Cost-Efficient Multi-Tenant Mobile Edge Computing

Shuaibing Lu¹, Member, IEEE, Jie Wu², Fellow, IEEE, Pengfan Lu, Ning Wang³, Haiming Liu⁴, and Juan Fang⁵, Member, IEEE

Abstract—The vigorous development of IoT technology has spawned a series of applications that are delay-sensitive or resource-intensive. Mobile edge computing is an emerging paradigm that provides services between end devices and traditional cloud data centers to users. However, with the continuously increasing investment of demands, it is nontrivial to maintain a higher quality-of-service (QoS) under the erratic activities of mobile users. In this paper, we investigate the service provisioning and updating problem under the multiple-users scenario by improving the performance of services with long-term cost constraints. We first decouple the original long-term optimization problem into a per-slot deterministic one by using Lyapunov optimization. Then, we propose two service updating decision strategies by considering the trajectory prediction conditions of users. Based on that, we design an online strategy by utilizing the committed horizon control method looking forward to multiple slots predictions. We prove the performance bound of our online strategy theoretically in terms of the trade-off between delay and cost. Extensive experiments demonstrate the superior performance of the proposed algorithm.

Index Terms—Cost-efficient, mobile edge computing, online service provisioning, quality -of-service (QoS).

I. INTRODUCTION

THE vigorous development of Internet of things (IoT) technology has led to the explosive growth of mobile terminal equipment and data volume. At the same time, a series of resource-intensive and delay-sensitive applications, such as augmented reality (AR)/virtual reality (VR), intelligent driving, and dynamic content delivery, have emerged and been widely used [1], [2], [4], [6]. It is difficult for the traditional cloud data

center to meet the performance requirements due to the long distance from massive terminals. Mobile Edge Computing (MEC) is a promising framework for solving this problem by deploying edge servers at base stations to supply computation, storage, and networking resources for multiple users [3]. Ensuring the quality of service (QoS) is a key challenge with significant real-world implications, as the vigorous development of IoT technology has generated numerous delay-sensitive or resource-intensive applications. However, the finite capabilities of edge servers and the erratic activities of multiple end-users pose challenges in guaranteeing the QoS. Therefore, there are two key problems: (i) How to guarantee the QoS to avoid service interruption with unknown trajectories when users are away from the original edge servers? (ii) How to realize service provisioning, and updating the services that can efficiently utilize the limited resources without overwhelming the cost constraint? In this paper, we investigate the service provisioning and updating problem under the multiple-users scenario by improving the performance of services with a long-term cost constraint.

A. Motivation and Challenges

Service provisioning and updating refers to the decision-making process of the locations of services within a specific edge computing network to balance the interests of service providers and consumers to achieve the greatest efficiency possible [7]. Numerous mobile devices and sensors in edge networks need to interface with services and exchange data in real-time, which requires efficient service provisioning and updating strategy to enhance the capacity to accommodate real-time data processing. We illustrate the motivation and challenges of the online service provisioning and updating problem by using an example in Fig. 1. The squares with six different colors represent the services s_1 to s_6 , which are initially provisioned in the cloud data center. We assume that the services required by the users have been deployed on the edge servers, and each service only serves one user. For mobile users, the QoS can be guaranteed through provisioning a replication or migration among edge servers. (i) The trajectories of multiple users are diverse and erratic, hence it is non-trivial to find an efficient strategy that can improve the QoS of mobile users by considering the cost constraint. Taking service s_3 as an example, we suppose that end user u_3 moves from an area in m_1 to m_4 at time slot t

Manuscript received 28 December 2022; revised 6 November 2023; accepted 19 November 2023. Date of publication 22 November 2023; date of current version 6 February 2024. This work was supported in part by the Fundamental Research Funds for the Central Universities under Grant 2021RC258, in part by the China Postdoctoral Science Foundation under Grant 2021M700366, and in part by the National Natural Science Foundation under Grant 92267107. (Corresponding author: Haiming Liu.)

Shuaibing Lu, Pengfan Lu, and Juan Fang are with the Faculty of Information Technology, Beijing University of Technology, Beijing 100124, China (e-mail: lshuaibing@bjut.edu.cn).

Jie Wu is with the Department of Computer and Information Sciences, Temple University, Philadelphia, PA 19122 USA (e-mail: jiewu@temple.edu).

Ning Wang is with the Computer Science at Rowan University, Glassboro, NJ 08028 USA (e-mail: wangn@rowan.edu).

Haiming Liu is with the School of Software Engineering, Beijing Jiaotong University, Beijing 100091, China (e-mail: liuhaiming@bjtu.edu.cn).

Digital Object Identifier 10.1109/TSC.2023.3335412

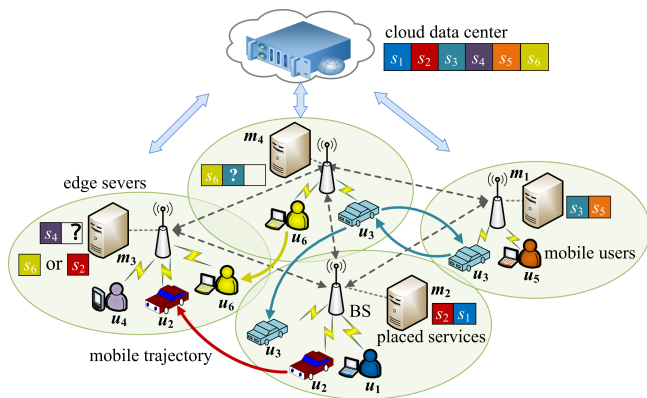


Fig. 1. Illustrating example.

and goes back to m_1 after several slots. One extreme solution is to migrate or provision a replication of s_3 on edge server m_4 which may bring a lower delay for user u_3 . However, the total cost will be the maximum one among all feasible assignments if the replication or migration costs of services are extremely high. Another extreme assignment is to retain service s_3 within m_1 , which minimizes the extra (replication or migration) cost. When u_3 moves to m_4 , the service can only be enjoyed through communication with m_1 , which will make the quality of service decrease. Therefore, when and where to migrate or replicate services is crucial for balancing the trade-off between long-term cost and users' total delay. (ii). Since the capabilities of edge servers are limited, determining which services are chosen to be placed in order to obtain a better performance when multiple users make the same decision at the same time is non-trivial. Taking services s_2 , s_3 , and s_6 as examples, we suppose that users u_2 and u_6 move from m_1 towards m_4 during the same time slot. If both services want to migrate or replicate to m_3 , there will be a conflict due to the fact that the remaining capacity can only receive one service. Let's consider a more serious scenario such that if m_3 finally agrees to allow service s_6 to migrate or replicate on itself, which means that service s_2 is still locating on m_2 , while during that time slot, if user u_3 moves from server m_4 to m_2 , there will be no choice for s_3 due to the limitation of capacity on m_2 . Therefore, the problem of how to make a better choice by jointly considering the resource efficiency and users' performance is a challenge.

B. Contributions and Paper Organization

In this paper, we investigate the online service provisioning and updating problem under the multiple-users scenario by improving the performance of services with long-term cost constraints in mobile edge computing. Our contributions can be summarized as follows:

- We investigate the online service provisioning and updating problem by formulating to minimize the average long-term delay of multiple mobile users, and we decouple the original long-term optimization problem into a per-slot deterministic one by using Lyapunov optimization.

- We propose two service updating decision strategies by considering the trajectory prediction conditions of users. For one scenario, namely the service updating without available predictive information, we propose a novel strategy by introducing the conflict resolution factor. For the other scenario, which is the service updating with multi-step prediction, we optimize the total delay of users per-slot by converting a weighted graph under the constructed activity set.
- Based on that, we design an online strategy by utilizing the committed horizon control method looking forward to multiple slots predictions. We prove the performance bound of our online service provisioning strategy theoretically in terms of the trade-off between delay and cost.
- We conduct extensive experiments to compare our strategy with several baselines based on the Microsoft GPS trajectory dataset which was reconstructed by 40 users. The results are shown from different perspectives to provide conclusions. Extensive experiments demonstrate the superior performance of the proposed algorithm.

The remainder of this paper is organized as follows. Section II surveys related works. Section III describes the model and then formulates the problem. Section IV investigates the service provisioning and updating problem based on Lyapunov optimization. Section V investigates the online optimization provisioning strategy by considering multi-step prediction. Section VI includes the experiments. Finally, Section VII concludes the paper.

II. RELATED WORK

As an emerging paradigm, edge computing extends services closer to end-users. However, the finite capabilities of edge servers and the erratic activities of users pose new challenges [5]. One of the main open branches is the service provisioning problem, which is well-investigated in edge computing under mobility scenarios [6]. Various works have been studied from different aspects of this problem. Elgazzar et al. [8] examined the reliability of mobile devices as data service providers and proposed a cloud-assisted framework that leverages task offloading to improve service performance. Qiu et al. [9] addressed the deployment optimization of VNFs and backup VNFs in a fault-prone MEC environment with dynamically changing fault probabilities. Li et al. [10] concentrated on the reliability of Virtual Network Functions (VNFs) in the MEC network and suggested deploying primary and backup VNF instances to meet user reliability requirements. Yu et al. [11] investigated the service provisioning problem in mobile edge computing, which aimed to minimize the traffic load caused by service request forwarding, and proposed an efficient decentralized algorithm based on matching theory. Mao et al. [12] proposed an approximate algorithm to deploy service function chains at the edge and the cloud, and they used the next fit strategy and double spanning tree to effectively avoid redundant data traffic and reduce the latency. The resource cost at the edge and on the cloud, and the optimal point of all corresponding communication delays were optimized. Gu et al. [13] proposed

the JMDLS-RR algorithm, which cooperatively deploys microservices by combining intra-server and inter-server layer sharing to maximize the service capacity of the edge cloud. Nezami et al. [14] formulated a decentralized load-balancing problem for IoT service provisioning, and they introduced a decentralized multi-agent system that utilized edge servers to balance the workloads and minimized the costs involved in service execution. Zhang et al. [15] solved the computation and delay costs minimization problem by proposing an efficiently approximate algorithm based on semi-definite relaxation. The above works optimized the cost and delay of services from the offline scenario.

In the online scenario, Chen et al. [16] studied the service collaboration with master-slave dependency among service chains of mobile users and jointly optimized the cost and delay by introducing a distributed algorithm based on Markov approximation. Xu et al. [17] proposed an efficient online algorithm based on Gibbs sampling which can achieve provable close-to-optimal performance. Ren et al. [18] proposed a novel framework called EdgeMatrix and designed a multi-task mechanism to solve the problem of joint service orchestration and request scheduling between edge cloud clusters. They used a dual timescale framework which coordinated resources and services on a large timescale and scheduled requests on a small one, which can significantly shorten the running time. Shang et al. [19] designed an online service provisioning and throughput adjustment algorithm to coordinate the migration of virtual services, as well as adjusted its data throughput according to real-time bandwidth fluctuations to reduce latency and improve the Quality of Experience (QoE). They solved the support challenge of interactive virtual service QoE in mobile edge computing caused by high user mobility and unstable network conditions. Wang et al. [20] described dynamic task placement as an online multi-user doobly slot machine process and proposed a decentralized algorithm to optimize users' rewards affected by network delays. Han et al. [21] transformed the online multi-component service placement into an ant colony optimization problem, and they proposed a level traversal component ranking method to achieve faster convergence. In the online optimization problem of edge computing, some existing works utilize the Lyapunov optimization method. Li et al. [22] proposed a two-timescale algorithm based on Lyapunov optimization, which achieves efficient performance close to offline optimal results by purchasing computing resources and making task offloading decisions at different timescales. Liu et al. [23] discussed the challenge of obtaining feasible computation offloading strategies in an online environment due to limited resources of unmanned aerial vehicles (UAVs) and dynamic changes in applications and environments, which achieved long-term efficient and stable performance. In cloud computing system optimization, the application of Lyapunov optimization techniques is also a promising approach. Zhou et al. [24] proposed an analytical framework for optimizing the trade-off between power consumption and performance in SaaS cloud platforms, and they used Lyapunov optimization techniques to design an optimal control framework for online decision-making on request admission control. Fang

et al. [25] designed a stochastic control algorithm using Lyapunov optimization and weight perturbation techniques, which achieved the maximization of profit for the management platform through online decision-making. Qi et al. [26] proposed an online scheduling algorithm based on Lyapunov optimization to optimize the operation of service systems in a cloud environment, utilizing queue stability to ensure Quality of Service (QoS). These works focus on optimizing the cost and delay of the service provisioning problem; however, they ignore the erratic movements of users.

In order to tackle the challenge of users' mobility, some existing works were proposed based on service migration. Ning et al. [27] studied the service provisioning problem by constructing a stochastic mobility system, and they introduced a distributed Markov approximation algorithm which is linear to the number of users in order to determine the configurations of services provisioning. Kim et al. [28] proposed a system called MoDEMS to optimize the service provisioning based on user mobility in edge computing, and they developed a linear integer programming problem and a Seq-Creedy heuristic method to generate a migration plan to minimize system costs and user delays. Zeng et al. [29] formulated an optimization problem to jointly decide the service provisioning policy and the routing decision, and they developed an online distributed algorithm with provable performance guarantees in terms of convergence and competitive ratio. Li et al. [30] focused on the service migration problem for mobile users through modeling a Markov Decision Process (MDP) model, and they solved it by using deep reinforcement learning. In addition, some works consider using the information of the prediction. Liu et al. [31], introduced a prediction-based dynamic task assignment algorithm that assigned the workloads to edge servers based on the prediction of capacities and costs in each time slot. Jin et al. [32] designed a set of novel polynomial-time algorithms to make adaptive decisions by solving continuous solutions. These continuous solutions are based on the predicted inputs about the dynamic and uncertain cloud-edge environments via online learning. Ma et al. [33] propose a multiple-slots predictive service placement algorithm to incorporate the prediction of user mobility based on a frame-based design. However, these works do not take into account the impact of additional prediction error on the service provisioning. In this paper, we study the online service provisioning and updating problem in mobile edge computing. Our objective is to improve the QoS by minimizing the total delay while considering maintaining the long-term cost under the constraint.

III. MODEL AND PROBLEM FORMULATION

In this paper, we focus on the QoS-aware online service provisioning and updating problem in mobile edge computing for multiple users with cost efficiency. We build a system model that describes the physical edge nodes and multiple users, and then we formalize our problem to minimize the long-term average delay under the constraints within this model.

TABLE I
SYMBOLS AND DEFINITIONS

Symbols	Definitions
\mathbf{M}	Set of mobile edge servers, where $\mathbf{M} = \{m_j\}$.
\mathbf{S}	Set of services on cloud data center, where $\mathbf{S} = \{s_h\}$.
$\mathbf{S}_{m_i}(t)$	Set of services provisioning on m_i at t .
\mathbf{U}	Set of mobile users, where $\mathbf{U} = \{u_i\}$.
$\hat{\mathbf{U}}(t)$	The activity set of users at time slot t .
$\mathbb{D}_{u_i}(t)$	Total delay of user u_i at time slot t .
$D_{u_i}^c(t)$	Computing delay of user u_i at time slot t .
$D_{u_i}^l(t)$	Communication delay of user u_i at time slot t .
$D_{u_i}^u(t)$	Updating delay of user u_i at time slot t .
$\mathbb{C}_{s_h}(t)$	Total cost of s_h at time slot t .
$C_{s_h}^m(t)$	Migration cost of s_h at time slot t .
$C_{s_h}^r(t)$	Replication cost of s_h at time slot t .
η_h	The conflict resolution factor of service s_h .
$L_{u_i}(t)$	The location of u_i at time slot t .
$\tilde{L}_{u_i} \tau, \tau+\omega$	The trajectory of u_i in a ω steps prediction window.

A. System Model

As shown in Fig. 1, we consider a three layer network architecture that includes the cloud data center, edge servers, and the mobile end-users. We suppose that the services required by users are initially provisioning in the cloud data center, which is denoted as set $\mathbf{S} = \{s_h\}$. Let $\mathbf{M} = \{m_j\}$ denote a substrate set of edge servers that are supported by the operators. Let $\mathbf{U} = \{u_i\}$ denote the set of mobile users, and these users subscribe to the services one-to-one. In order to capture the mobility of users, we assume that the system in a slotted structure and its timeline is discretized into time frame $t \in \{0, 1, 2, \dots, T-1\}$ [33], [34], [37]. In this paper, we suppose that users move erratically and frequently among several edge servers. At each time slot, the operators determine whether provisioning replications or migration follow with users according to navigating the trade-off between delay and cost. We list the main notations throughout this paper in Table I for ease of reference.

1) *QoS Model*: In our study, the QoS of users is determined by computing delay, communication delay, and updating delay. We use $\mathbb{D}(t) = \sum_{i=1}^{|\mathbf{U}|} \mathbb{D}_{u_i}(t)$ to denote the total delay at time slot t , where $\mathbb{D}_{u_i}(t)$ is the delay of u_i . The computing delay is defined as $D_{u_i}^c(t) = \sum_{m_j \in \mathbf{M}} \frac{r_{u_i}(t)}{z_{m_j}^c}$, where $r_{u_i}(t)$ is the service request of user u_i at time slot t , and $z_{m_j}^c$ is the computing capacity of m_j measured by the number of CPU cycles [35], [36]. We use $D_{u_i}^l(t)$ to represent the communication delay that occurs when users are far away from the location of the service. Let t_{u_i, m_j} denote the maximum transmission rate between user u_i and edge server m_j , where $t_{u_i, m_j}(t) = b_{u_i, m_j}(t) \cdot \log_2(1 + \frac{\beta \cdot g(u_i, m_j)}{N})$ [37], [38], [44], [45]. We set $b_{u_i, m_j}(t)$ as a binary indicator variable indicating whether user u_i is connected to server m_j . Here, by β we denote the transmission power of the local mobile device of u_i . Let $g(u_i, m_j)$ be the channel gain between the user u_i and the edge server m_j , where $g(u_i, m_j) = 127 + 30 \cdot \log p(u_i, m_j)$ [39]. $p(u_i, m_j)$ represents the distance between u_i and m_j that determine the network propagation. We use N to represent the noise power. The communication delay is defined as $D_{u_i}^l(t) = \sum_{m_j \in \mathbf{M}} \frac{d_{u_i}(t)}{t_{u_i, m_j}(t)}$, where $d_{u_i}(t)$ denotes the data size of the request [37], [38]. We use $D_{u_i}^u(t)$ to represent the updating delay, which occurs when the location of service s_i

that is serving u_i changes. Here, we consider two scenarios. One is that the operator can place a replication on the edge server to which u_i is currently connected. The other is that operator can migrate service s_i to the edge server to which user u_i goes forward. The costs of both scenarios are discussed in the next subsection. The updating delay is defined as $D_{u_i}^u(t) = \Upsilon(v_i) + \Psi(s_i)$, where $\Upsilon(s_i)$ is the delay of rebooting software resources, and $\Psi(s_i)$ is the delay of transmitting service profiles [39].

2) *Cost Model*: We use $\mathbb{C}(t)$ to denote the total cost of users in set \mathbf{U} at time slot t , where $\mathbb{C}(t) = \sum_{h=1}^{|\mathbf{S}|} \mathbb{C}_{s_h}(t)$. Let $\mathbb{C}_{s_h}(t)$ denote the cost of service s_h , where $\mathbb{C}_{s_h}(t) = C_{s_h}^m(t) + C_{s_h}^r(t)$. We use $C_{s_h}^m(t)$ and $C_{s_h}^r(t)$ to represent the migration cost and replication cost, respectively. Let $x_{s_h}(t)$ denote the decision of s_h , when s_h decides to stay at the edge server with the same location in the previous step, $x_{s_h}(t) = 0$, otherwise, $x_{s_h}(t) = 1$.

B. Problem Formulation

On the basis of the models above, our problem is formulated to minimize the long-term average delay subject to the resource and cost constraints, which is presented as follows:

$$\mathbf{P}_1 : \text{minimize } \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \sum_{i=1}^{|\mathbf{U}|} \mathbb{D}_{u_i}(t) \quad (1)$$

$$\text{s.t. } \mathbb{D}_{u_i}(t) = D_{u_i}^c(t) + D_{u_i}^l(t) + x_{s_h}(t) \cdot D_{u_i}^u(t), \quad (2)$$

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^T \sum_{h=1}^{|\mathbf{S}|} \mathbb{C}_{s_h}(t) \leq \bar{\Gamma}, \mathbb{D}_{u_i}(t) \leq \bar{D}, \forall u_i \in \mathbf{U}, \quad (3)$$

$$\sum_{\mathbf{S}_{m_i} \in \mathbf{S}} W(\mathbf{S}_{m_i}(t)) \leq R_{m_i}^s, \forall m_i \in \mathbf{M}, \quad (4)$$

$$x_{s_h}(t) \in \{0, 1\}, \forall s_h \in \mathbf{S}. \quad (5)$$

\mathbf{P}_1 is the objective function, and equations (2) to (5) are the constraints. Equation (2) is the total delay of each user, which needs to be lower than \bar{D} to ensure the QoS. Equation (3) states that the long-term average cost cannot exceed the threshold $\bar{\Gamma}$ determined by the operators. Equation (4) states the constraint on the resource, which means the services placed on m_i should be under the limitation $R_{m_i}^s$. Here, we use $W(\mathbf{S}_{m_i}(t))$ to denote the amount of storage resources occupied for provisioned services on edge server m_i . Equation (5) states the decision of s_h which provides service for u_i at time slot t . As shown in the above equations, in order to obtain the optimal solution of \mathbf{P}_1 , complete offline information is required, i.e., the distribution of users' trajectories over all time slots, which is difficult to predict in advance. Thus, the main challenge that complicates the derivation of the optimal solution to the above problem is the lack of future information. In addition, the constraints that \mathbf{P}_1 must satisfy during the optimization process that make it very difficult to solve even if the future information is known a priori. Therefore, these challenges require an online approach that enables service provisioning and updating decisions to be made efficiently.

IV. SERVICE UPDATING DECISION STRATEGY BASED ON LYAPUNOV OPTIMIZATION

In this section, we first introduce two service updating decision strategies based on Lyapunov optimization by considering the trajectory prediction condition of multiple mobile users.

A. Decoupling Based on Lyapunov Optimization

Since the major challenge of directly solving \mathbf{P}_1 is that the long-term cost constraint of providers couples the service provisioning and updating decisions across different time slots, we first decouple the original problem into per-frame deterministic problems by applying Lyapunov optimization in this subsection. In order to deal with the constraint on average cost \bar{C} in (3), we introduce a virtual queue $Q(t)$ which denotes the historical measurement of the extra cost of services at time slot t . The queue updates according to the following equation

$$Q(t+1) = \max\{Q(t) + C(t) - \bar{C}, 0\} \quad (6)$$

Intuitively, the condition of the total extra cost $C(t)$ that is produced by the replication or migration of services can be evaluated by $Q(t)$. When the value of $Q(t)$ is large, it represents that the cost has exceeded the long-term cost \bar{C} . Specifically, (6) implies $Q(t+1) \geq Q(t) + C(t) - \bar{C}$, and then we have $C(t) - \bar{C} \leq Q(t+1) - Q(t)$. By summing this inequality during all time slots, we have $\sum_{t=0}^{T-1} (C(t) - \bar{C}) \leq Q(T) - Q(0)$. Initialize $Q(0) = 0$ and divide by T time slots. One can take expectations and derive that the expected backlog over time slot in $[0, T-1]$ is less than the threshold.

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[C(t)] \leq \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}[Q(T)] + \bar{C} \quad (7)$$

As shown in (7), we have that the constraint on the cost can be guaranteed by stabilizing the virtual queue $Q(t)$. Therefore, a quadratic Lyapunov function for each slot t is defined as $L(Q(t)) \triangleq \frac{1}{2}Q(t)^2$ [27], [33], [40], where $Q(t)$ is a vector that evolves over slots in $[0, T-1]$. Here, the quadratic Lyapunov function can be considered as a scalar measure of queue deviation which is similar to $Q(t)$. In order to keep the queue stable, which means enforcing the extra cost constraint by promoting the Lyapunov function to lower values continuously, we introduce the one-step conditional Lyapunov drift as follows.

$$\Delta(Q(t)) \triangleq \mathbb{E}[L(Q(t+1)) - L(Q(t))|Q(t)] \quad (8)$$

Lemma 1: Given the updating decisions of services in set \mathbf{S} according to multiple mobile users \mathbf{U} in each time slot t , the statement holds:

$$\Delta(Q(t)) \leq \beta + Q(t)\mathbb{E}[(C(t) - \bar{C})|Q(t)] \quad (9)$$

where $\beta \triangleq \frac{1}{2}(\bar{C}^2 + \bar{C}^2)$.

Proof: We rearrange (8) for a concise form, where $\Delta(Q(t)) \triangleq \mathbb{E}[L(Q(t+1)) - L(Q(t))|Q(t)] = \frac{1}{2}\mathbb{E}[(C(t) - \bar{C})^2|Q(t)] + Q(t)\mathbb{E}[(C(t) - \bar{C})|Q(t)]$. For each service, we use $\tilde{C}_{s_h}(t)$ to denote the cost of updating the decision of s_h in set \mathbf{S} by choosing the minimum delay of user $u_h \in \mathbf{U}$ at time slot t . Based on that, the total cost of all services

will be $\tilde{C}(t) = \sum_{s_h \in \mathbf{S}} \tilde{C}_{s_h}(t)$. Because of the division of the time space taking into account the user's mobility on the boundary, the service provider will not change in one time slot. Thus, we have $\tilde{C}(t) \geq C(t)$. Then, we have $\Delta(Q(t)) \leq \frac{1}{2}(\tilde{C}(t) - \bar{C})^2 + Q(t)\mathbb{E}[(C(t) - \bar{C})|Q(t)] \leq \frac{1}{2}(\tilde{C}(t)^2 + \bar{C}^2) + Q(t)\mathbb{E}[(C(t) - \bar{C})|Q(t)]$. Therefore, we can obtain that the one-step conditional Lyapunov drift holds $\Delta(Q(t)) \leq \beta + Q(t)\mathbb{E}[(C(t) - \bar{C})|Q(t)]$ at each time slot t , where $\beta \triangleq \frac{1}{2}(\tilde{C}(t)^2 + \bar{C}^2)$. Therefore, the proof of Lemma 1 is complete. ■

According to the Lyapunov optimization framework, we obtain the upper bound of the Lyapunov drift function by introducing a Lyapunov drift-plus-penalty function in each time slot t .

$$P(t) \triangleq \Delta(Q(t)) + V\mathbb{E}[\mathbb{D}(t)|Q(t)] \quad (10)$$

Here, we define V as a non-negative parameter for adjusting the trade-off between the extra cost queue and the delay. In each time slot, the performance of the service provisioning strategy is guaranteed by minimizing an upper bound of the following function.

$$P(t) \leq \beta + Q(t)\mathbb{E}[(C(t) - \bar{C})|Q(t)] + V\mathbb{E}[\mathbb{D}(t)|Q(t)] \quad (11)$$

Based on that, the service provisioning and updating problem is formulated by minimizing the right side of (11) at each time slot, which is formulated as follows.

$$\mathbf{P}_2 : \text{minimize } \beta + Q(t)(C(t) - \bar{C}) + V\mathbb{D}(t) \quad (12)$$

$$\text{s.t. (2)–(5).} \quad (13)$$

B. Optimal Services Updating Decision Strategy

In this subsection, we propose a service updating decision strategy by optimizing \mathbf{P}_2 under the constraints in each time slot. We start with a definition as follows.

Definition 1 (Optimal Service Updating (OSU) Problem): Given the distribution of users \mathbf{U} , the topology of edge network \mathbf{G} , and the function $\Theta(t)$, an OSU problem is how to find a decision for services in \mathbf{S} to minimize \mathbf{P}_2 under the constraints at time slot t .

On the basis of Definition 1, we discuss two scenarios. One is the services updating without prediction, and the other is the service updating with prediction.

1) *OSU With No Prediction:* The first scenario we considered is the OSU problem without available information caused either by the inaccurate prediction results, or by it being the initial or training stages of mobile users in per-slot. The specific steps are shown in Algorithm 1. We use the sets of edge servers \mathbf{M} , users \mathbf{U} , and services \mathbf{S} as the input. The output is the service updating decision $\mathbf{X}(t)$ at time slot t . For each user in set \mathbf{U} , we choose the updating decision by optimizing \mathbf{P}_2 in lines 1 to 2. Then, we check the feasibility of services on edge servers by checking whether $\sum_{s_{m_i} \in \mathbf{S}} W(s_{m_i}(t)) \geq R_{m_i}^s$. Here, we use $\sum_{s_{m_i} \in \mathbf{S}} W(s_{m_i}(t))$ to denote the total number of services provisioning on m_i . In order to avoid conflicts caused by aggregation requests of multiple users, we introduce a definition

Algorithm 1: Updating Strategy with No Prediction (USNP).

Input: Sets of edge servers \mathbf{M} , users \mathbf{U} , and services \mathbf{S} ;
Output: Service updating decision $\mathbf{X}(t)$ of \mathbf{U} at time slot t ;

- 1: **for** users $k = 1$ to $k = |\mathbf{U}|$ in \mathbf{U} **do**
- 2: Choose the updating decision by optimizing \mathbf{P}_2 ;
- 3: **for** edge servers $i = 0$ to $i = |\mathbf{M}|$ in \mathbf{M} **do**
- 4: **if** $\sum_{s_{m_i} \in \mathbf{S}} W(\mathbf{S}_{m_i}(t)) \geq R_{m_i}^s$ **then**
- 5: Choose service by $i = \arg \min\{\eta_h\}$;
- 6: **end if**
- 7: **end for**
- 8: **end for**
- 9: **return** Service updating decision $\mathbf{X}(t)$ of \mathbf{S} ;

of the conflict resolution factor for the service, and the specific definition is as follows.

Definition 2 (Conflict Resolution Factor): Let η_h indicate the conflict resolution factor of service s_h and $\eta_h = \mathbb{C}_{s_h}(t) / \overline{\mathbb{D}_{u_h}^l(t)}$, where $\overline{\mathbb{D}_{u_h}^l(t)} = \mathbb{D}_{u_h}^l(t) |_{s_h \notin \mathbf{S}_{m_i}(t)}$.

Here, we use $\mathbb{C}_{s_h}(t)$ to denote the total extra cost of service s_h when it migrates or replicates on edge server m_i at time slot t , where $s_h \in \mathbf{S}_{m_i}(t)$. In line 4, we choose a service by an increasing order $i = \arg \min\{\eta_h\}$. Finally, the service updating decision $\mathbf{X}(t)$ is returned in line 6.

2) *OSU With Prediction:* In this subsection, we explore a more realistic and complicated scenario in which we consider the trajectory prediction for the service provisioning and updating decision strategy. Here, we introduce an online strategy in the view of the committed horizon control method, where the predictions are looking forward to several slots for multiple slots by utilizing the existing well-performance model.

Lemma 2: The decision of OSU problem can be solved by minimizing $\Theta(t)$, where $\Theta(t) = Q(t)\mathbb{C}(t) + V\mathbb{D}(t)$.

Proof: We first rearrange \mathbf{P}_2 by introducing an intermediate variable \mathbb{P} , where $\mathbb{P}(t) = \beta + Q(t)(\mathbb{C}(t) - \bar{\Gamma}) + V\mathbb{D}(t) = \beta + Q(t)\mathbb{C}(t) - Q(t)\bar{\Gamma} + V\mathbb{D}(t)$. The value of β is related to the distribution of users in set \mathbf{U} , which is a constant value. Meanwhile, the value of $Q(t)$ depends on the decision of services in the previous time interval $[0, t-1]$, which means that the decision at time slot t has no effect on the value of $Q(t)$. We reconstruct $\mathbb{P}(t)$ as $\mathbb{P}(t) = W + \Theta(t)$, where $W = \beta + Q(t) - \bar{\Gamma}$ and $\Theta(t) = Q(t)\mathbb{C}(t) + V\mathbb{D}(t)$. Therefore, we can obtain that the network determines the service updating strategies by solving the optimization of $\Theta(t)$ in each time slot. ■

Based on the conversion above, we rearrange $\Theta(t)$ by considering the combinational decision-making where $\Theta(t) = Q(t) \sum_{h=1}^{|\mathbf{S}|} \mathbb{C}_{s_h}(t) + V \sum_{i=1}^{|\mathbf{U}|} \mathbb{D}_{u_i}(t)$. The value of the total extra cost of service s_h depends on the decision choosing to migrate or place replications, i.e., $\mathbb{C}_{s_h}(t) = C_{s_h}^m(t) + C_{s_h}^r(t)$, which will affect the result of the delay. Taking the decision of s_h as an example, if service s_h decides to migrate or place replications on other edge servers, it will produce a migration cost $C_{s_h}^m(t)$ or replication cost $C_{s_h}^r(t)$. Meanwhile, the communication part of $\mathbb{D}_{u_i}^l(t)$ will decrease while the updating part of $\mathbb{D}_{u_i}^u(t)$ will increase for $\mathbb{D}_{u_i}(t)$. We reconstruct $\Theta(t)$ on the basis

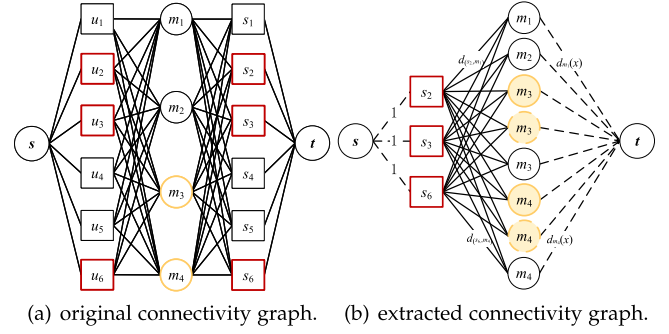


Fig. 2. Connectivity graphs of Fig. 1.

of the interaction based on the relationship between services and users, where $\Theta(t) = \sum_{h=1}^{|\mathbf{S}|} \Theta_h(t)$. For each service, we have $\Theta_h(t) = Q(t)\mathbb{C}_{s_h}(t) + (D_{u_h}^l + D_{u_h}^u) + D_{u_h}^c$.

Based on that, we use $d_{(s_h, m_i)}$ to represent the weight between service s_h and edge server m_i at time slot t , where $d_{(s_h, m_i)}(t) = Q(t)\mathbb{C}_{s_h}(t)$. We suppose that $d_{m_i}(x)$ is the delay function. We replace each edge in \mathbf{G}° with $|\hat{\mathbf{U}}(t)|$ parallel edges between the same server m_i and the destination t , and each with weight $d_{m_i}(x) |_{u_x \in \hat{\mathbf{U}}(t)}$. Then, the weight between edge server m_i and the destination t that is connected to it is $d_{m_i}(x) = (D_{u_x}^l + D_{u_x}^u) + D_{u_x}^c$. Therefore, we have $\Theta_h(t) = d_{(s_h, m_i)}(t) + d_{m_i}(x)(t)$.

On the basis of the interaction, we propose a novel Updating Strategy with No Prediction (USNP) to optimize the provisioning strategy at each time slot, which is shown in Algorithm 2. We first construct an original connectivity graph by considering the information and connection between services and edge servers. We add two virtual nodes which are source s and destination t , and the middle two layers are services and storage resources of edge servers. The original connectivity graph is shown in Fig. 2(a), where users are represented by the left squares in this diagram, edge servers by the middle circles, and services by the right squares. We assume that edge nodes are bidirectionally reachable, which means that users can access all edge nodes with different costs and delays. In addition, the service is able to choose any edge server for provisioning when the capacity of the edge server allows, however, different positions will result in different delays. Thus, the users and services are fully connected to the edge servers. In each time slot, the activities of users are dynamic and independent. This means that some users may be remaining in their original locations, while others may be far away from the connected edge servers. We use thick red lines to mark services where their users are far away from the original locations, and thick yellow lines to mark edge servers with remaining resources. For the users whose locations are not changing, the corresponding service will not be migrated or placed by an additional replica, so there is no extra cost or delay produced. Therefore, we consider optimizing the provisioning of services by constructing an activity set $\hat{\mathbf{U}}(t)$ to reduce the dimensional space. The formal definition is given as follows.

Definition 3 (Activity Set): Let $\hat{\mathbf{U}}(t)$ indicate the activity set of users at time slot t , where $u_i \in \hat{\mathbf{U}}(t)$ is the user whose current

Algorithm 2: Updating Strategy with Prediction (USP).

Input: Sets of edge servers \mathbf{M} , users \mathbf{U} , and services \mathbf{S} ;
Output: Service updating decision $\mathbf{X}(t)$ of \mathbf{S} at time slot t ;

- 1: Construct the original connectivity graph \mathbf{g} based on the provisioning of \mathbf{S} , the connections of \mathbf{G} , and \mathbf{U} ;
- 2: **for** users $i = 1$ to $i = |\mathbf{U}|$ in \mathbf{U} **do**
- 3: Calculate $\varsigma_{u_i}(t) = (L_{u_i}(t-1), L_{u_i}(t))$;
- 4: **if** $\varsigma_{u_i}(t) == 1$ **then**
- 5: Construct the activity set with $\hat{\mathbf{U}}(t) \leftarrow u_i$;
- 6: Update user set at time slot t with $\mathbf{U}(t) = \mathbf{U}(t)/u_i$;
- 7: **else**
- 8: Update $\mathbf{U}(t) \leftarrow u_i$;
- 9: **end if**
- 10: **end for**
- 11: Construct the extracted connectivity graph \mathbf{G}° based on the activity set $\hat{\mathbf{U}}(t)$;
- 12: Replace the link with $|\hat{\mathbf{U}}(t)|$ parallel ones with weight $d_{m_i}(x)|_{u_x \in \hat{\mathbf{U}}(t)}$;
- 13: Find a feasible service updating decision with min-cost flow of $\hat{\mathbf{U}}(t)$;
- 14: **return** Service updating decision $\mathbf{X}(t)$ of services \mathbf{S} ;

location $L_{u_i}(t)$ is going far away from the edge server for initial connection $L_{u_i}(t-1)$.

Here, we use $L_{u_i}(t)$ to denote the edge server that user u_i becomes connected to at time slot t . Since one user can only be served by one service, the number of users and services are equal. Based on that, we do an extraction by considering the current status of users and the topology of the edge network. The extracted connectivity graph is shown in Fig. 2(b). We use the white circle to indicate that a container on the edge server has been occupied. Due to the fact that each server provisioning service has two alternatives, replication and migration, these options are accessible. The yellow dashed circle indicates migration, and the edge node, where a copy is placed, is represented by a solid yellow circle. Based on that, we assign values to network edges. The virtual source node which is connected to all user services with the weight of 1. This means that only one service provisioning decision can be made for each service, i.e., only one server can be selected for one service. Here, the weight of a service and each edge node decision edge for migration is the sum of migration delay and cost, while the weight of replication decisions is the sum of replication cost and delay. Each edge server involves two service provisioning decisions, each of which is connected to a virtual destination node. Here, the weight of the migration decision between the edge server and the virtual destination is the migration delay, while the replication decision is the replication delay.

The specific steps are shown in Algorithm 2. We use the sets of \mathbf{M} , \mathbf{U} , and \mathbf{S} as the inputs. The service updating decision $\mathbf{X}(t)$ of \mathbf{S} at time slot t is used as the output. We construct the original connectivity graph \mathbf{g} based on the provisioning of \mathbf{S} , the connections of \mathbf{G} , and \mathbf{U} in line 1. In line 2, we start to construct

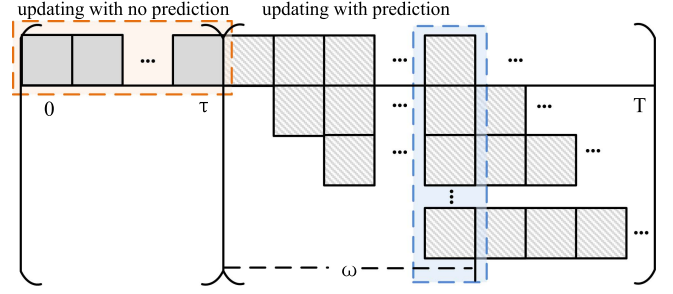


Fig. 3. Illustrating example of Algorithm 3.

the activity set $\hat{\mathbf{U}}(t)$. We first check the locations of users in set \mathbf{U} , where $\varsigma_{u_i}(t) = (L_{u_i}(t-1), L_{u_i}(t))$ in line 3. If $\varsigma_{u_i}(t) = 1$, this denotes that u_i has gone away from the edge server at time slot $t-1$. Then, we construct the activity set by adding u_i into set $\hat{\mathbf{U}}(t)$, where $\hat{\mathbf{U}}(t) \leftarrow u_i$; Otherwise, it denotes that u_i always stays near the edge server from $t-1$ to t , and we update $\mathbf{U}(t) \leftarrow u_i$. Based on this, we start to construct the extracted connectivity graph \mathbf{G}° based on the activity set $\hat{\mathbf{U}}(t)$ in line 9. In line 10, we replace the link with $|\hat{\mathbf{U}}(t)|$ parallel ones with weight $d_{m_i}(x)|_{u_x \in \hat{\mathbf{U}}(t)}$ between edge servers and destination t . Then, we find a feasible service updating decision with min-cost flow of $\hat{\mathbf{U}}(t)$, and we return the updating decision $\mathbf{X}(t)$ of services \mathbf{S} in line 12.

V. ONLINE OPTIMIZATION OF SERVICE PROVISIONING STRATEGY

In this section, we design an Online Optimization of Service Provisioning Strategy (O-OSP $_\omega$) by utilizing the committed horizon control method with ω steps prediction. The main ideas of O-OSP $_\omega$ are to leverage the prediction model to look forward the trajectories of users in multiple steps and to use the information to realize the optimization of service provisioning.

The specific steps are shown in Algorithm 3. We illustrate the whole process through Fig. 3. We suppose that the chosen prediction model exists a small part of adjustment stage in the initial τ time slots (orange covered area in Fig. 3), which means that the information in first τ steps is unavailable. Thus, we get service updating decision $\mathbf{X}(t)$ using Algorithm 1 in line 2. After that, we obtain the service updating decision $\mathbf{X}(t)$ using Algorithm 2 based on $\hat{\mathbf{L}}_{\mathbf{U}}|_{[t, t+\omega]}$ in line 4. Here, $\hat{\mathbf{L}}_{u_i}|_{[\tau, \tau+\omega]}$ is the trajectory of user u_i in a ω time steps prediction window starting at time τ , where $\hat{\mathbf{L}}_{u_i}|_{[\tau, \tau+\omega]} = \{\hat{L}_{u_i}(\tau), \hat{L}_{u_i}(\tau+1), \dots, \hat{L}_{u_i}(\tau+\omega)\}$ (blue covered area in Fig. 3). In line 5, we set $\tilde{t} = (t - \tau) \bmod \omega$, and we check whether the prediction steps are less than ω . In lines 9 to 13, we update the service provisioning for services by introducing a novel factor feasible decision frequency. We use $a_{s_h}(t) = x_{s_h}(t) \cdot y_{s_h}(t)$ to represent the decision value of s_h , where $y_{s_h}(t) \in \{-1, 1\}$ and $x_{s_h}(t) \in \{0, 1\}$ as shown in equation (5). Since $x_{s_h}(t) = 0$ when service s_h decides to stay at the original location, the decision value will be $a_{s_h}(t) = 0$. On the contrary, when service s_h decides on migration, $y_{s_h}(t) = -1$ and $x_{s_h}(t) = 1$, and then the value of $a_{s_h}(t) = -1$. Similarly,

Algorithm 3: Online Optimization of Service Provisioning strategy (O-OSP_ω).

Input: Sets of edge servers \mathbf{M} , users \mathbf{U} , and services \mathbf{S} ;

Output: Service updating decision \mathbf{X} of \mathbf{S} in each time slot;

```

1: for  $t = 0$  to  $t = \tau$  do
2:   Get service updating decision  $\mathbf{X}(t)$  using
   Algorithm 1;
3: end for
4: for  $t = \tau$  to  $t = T - 1$  do
5:   Get service updating decision  $\mathbf{X}(t)$  using
   Algorithm 2 based on  $\hat{\mathbf{L}}_{\mathbf{U}|[t, t+\omega]}$ ;
6:   Set  $\tilde{t} = (t - \tau) \bmod \omega$ ;
7:   if  $\tilde{t} = t - \tau$  then
8:     Set  $\mathbf{X}(t) = \mathbf{X}(\tilde{t})$ ;
9:   else
10:    for service  $h = 1$  to  $h = |\mathbf{S}|$  do
11:      Update the decision value of  $s_h$  into  $A_{s_h}^{(\omega)}$ ;
12:      Calculate the decision policy frequencies;
13:      Set  $X_{s_h}(\tilde{t}) = \arg \max_{a^\circ \in A_{s_h}^{(\omega)}} \{\varrho_{s_h|a^\circ}^{a^\circ}\}$ ;
14:    end for
15:    Set  $\mathbf{X}(t) = \{X_{s_h}(\tilde{t})\}_{s_h \in \mathbf{S}}$ ;
16:  end if
17: end for
18: return Service updating decision  $\mathbf{X}(t)$  of  $\mathbf{S}(t)$ ;

```

when service s_h makes a decision on replication, $y_{s_h}(t) = 1$ and $x_{s_h}(t) = 1$, and then the value of $a_{s_h}(t) = 1$. Based on that, we use a queue $A_{s_h}^{(x)}$ to record the decision values of service s_h in x time steps, i.e., $A_{s_h}^{(\omega)} = \{a_{s_h}(t+1), a_{s_h}(t+2), \dots, a_{s_h}(t+\omega)\}$.

Definition 4 (feasible decision frequency): Let $\varrho_{s_h|a^\circ}^{a^\circ}(t)$ indicate the feasible decision frequency of s_h under the value a° , where $\varrho_{s_h|a^\circ}^{a^\circ}(t) = \frac{1}{\omega} \sum_{x=0}^{x=\omega-1} f(A_{s_h}^{(x)}, a^\circ)$.

Here, $f(A_{s_h}^{(x)}, a^\circ)$ is a function to indicate whether the result in queue $A_{s_h}^{(x)}$ is equal to a° , i.e., $a_{s_h} = a^\circ$. Then, we choose the updating decision of s_h by setting $X_{s_h}(\tilde{t}) = \arg \max_{a^\circ \in A_{s_h}^{(\omega)}} \{\varrho_{s_h|a^\circ}^{a^\circ}\}$.

Theorem 1: By applying O-OSP_ω, the time-average system delay satisfies:

$$\frac{1}{T} \sum_{t=0}^{t=T-1} \mathbb{D}(t) \leq \frac{1}{2} (OPT + \beta + V|\mathbf{U}|\bar{D}) + \epsilon + \frac{1}{\omega} W \cdot \alpha \cdot T.$$

Proof: We conduct the proof via introducing $\mathbb{P}_{OSP}(t)$, where $\mathbb{P}_{OSP}(t) = \beta + Q(t)(\mathbb{C}(t) - \bar{\Gamma}) + V\mathbb{D}(t)$ under the O-OSP_ω strategy. For each time slot, we use $\mathbb{P}(t)$ to represent the decision policy with random frequency. We use $\delta(t)$ to denote the prediction error at time slot t . Then, we have the average value $\frac{1}{\omega} \sum_{t+1}^{t+\omega} b(t) \leq \frac{1}{\omega} \cdot \omega \cdot \epsilon = \epsilon$. Thus, we have

$$\mathbb{P}_{OSP}(t) \leq \frac{1}{\omega} \sum_{t+1}^{t+\omega} \mathbb{P}(t) \leq OPT + 2\epsilon + \frac{2}{\omega} W \cdot \alpha \cdot T, \quad (14)$$

which can be obtained in [41]. Here, $W = \max_{m_i \in \mathbf{M}} \{W(\mathbf{S}_{m_i}(t))\}$, which denotes the maximum available storage resource of edge servers. Since each server needs to make decisions for the services that are placed on it, there exists a stationary and randomized policy π for \mathbf{P}_2 which satisfy $\mathbb{E}[\mathbb{C} - \bar{\Gamma}] \leq \delta$. Thus, we have $\mathbb{P}_{OSP}(t) \leq \beta + Q(t) \cdot \delta + V\mathbb{D}(t)$. By letting δ go to zero, we have $\mathbb{P}_{OSP}(t) \leq \beta + V\mathbb{D}(t)$.

$$\mathbb{P}_{OSP}(t) \leq \beta + V\mathbb{D}(t) \leq \beta + V|\mathbf{U}|\bar{D}. \quad (15)$$

We sum the inequalities labeled as Equations (14) and (15), and then we have

$$\mathbb{P}_{OSP}(t) \leq \frac{1}{2} (OPT + \beta + V|\mathbf{U}|\bar{D}) + \epsilon + \frac{1}{\omega} W \cdot \alpha \cdot T. \quad (16)$$

Therefore, the proof of Theorem 1 is complete. \blacksquare

In this paper, the O-OSP_ω algorithm uses ω steps trajectory prediction and min-cost flow; the time complexity will be determined by the maximum value of these two parts. Since the trajectories of users are predicted through social-LSTM, and the provisioning process is based on the predicted results. The time complexity is $O(\mathbf{T}\mathbf{B}\mathbf{H}^2)$, where \mathbf{T} is the length of the sequence determined by the steps ω , i.e., $\mathbf{T} = 2 + 2\omega$, \mathbf{B} is the batch size, and \mathbf{H} is the scale of the network hidden layer. In addition, the complexity of min-cost flow is $O(\mathbb{G}\mathbb{E}\mathcal{T})$, where \mathbb{G} is the total number of nodes in the extracted connectivity graph, i.e., $\mathbb{G} = 2 + |\hat{\mathbf{U}}(t)| + |\bar{\mathbf{M}}| + \sum_{m_i \in \mathbf{M}/\bar{\mathbf{M}}} R_{m_i}^{s_i}$. Here, $|\bar{\mathbf{M}}|$ is the total number of edge servers which are fully occupied. \mathbb{E} is the total number of edges in the extracted connectivity graph, i.e., $\mathbb{E} = |\hat{\mathbf{U}}(t)| + |\hat{\mathbf{U}}(t)|(|\bar{\mathbf{M}}| + \sum_{m_i \in \mathbf{M}/\bar{\mathbf{M}}} R_{m_i}^{s_i}) + |\bar{\mathbf{M}}| + \sum_{m_i \in \mathbf{M}/\bar{\mathbf{M}}} R_{m_i}^{s_i} = (|\hat{\mathbf{U}}(t)| + 1)(|\bar{\mathbf{M}}| + \sum_{m_i \in \mathbf{M}/\bar{\mathbf{M}}} R_{m_i}^{s_i} + 1) - 1$. \mathcal{T} is the number of times the cumulative augmented path reaches the maximum flow. Therefore, the complexity of O-OSP_ω is $O(\max\{\mathbf{T}\mathbf{B}\mathbf{H}^2, \mathbb{G}\mathbb{E}\mathcal{T}\})$, which is the maximum value of both trajectory prediction and min-cost flow.

VI. EXPERIMENTS

In this section, we conduct the experiments based on the Microsoft GPS trajectory dataset [42], [43] to study the service provisioning problem for multiple mobile users in edge computing networks.

A. Basic Setting

We build our prototype on a workstation that runs a Linux operating system with E5-2620 CPU, NVIDIA RTX5000 GPU, 128 Gb memory, and a 2 Tb hard disk. We choose the Social-LSTM model to predict the future trajectories of users which can achieve an average accuracy of over 70%. We used the published Microsoft GPS trajectory dataset which has been collected in the Geolife project [42], [43]. The Microsoft GPS trajectory dataset is a GPS trajectory dataset collected from 182 users over a period of more than three years (from April 2007 to August 2012) as part of the (Microsoft Research Asia) Geolife project. These trajectories were recorded by various GPS loggers and GPS phones, and they have a variety of sampling rates.

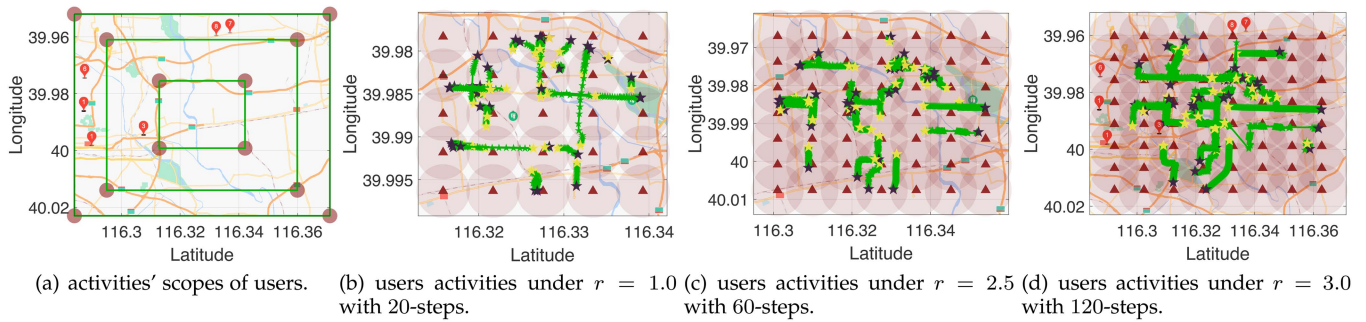


Fig. 4. Users' activities under different scopes of users with scaling steps.

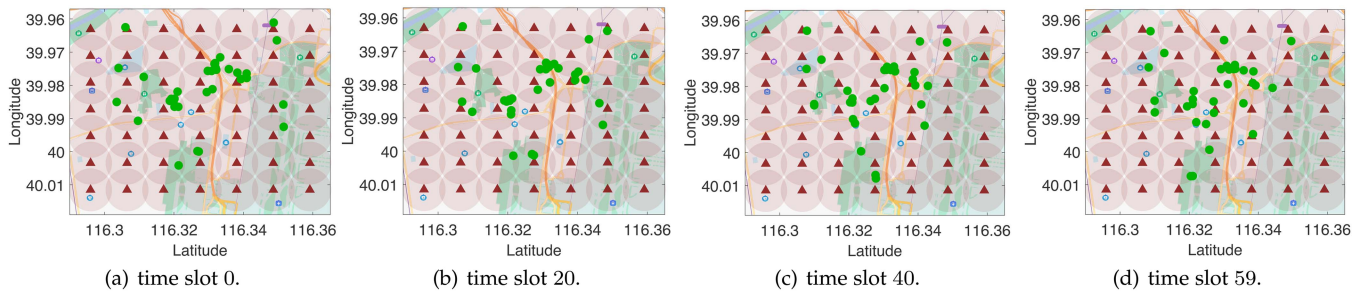


Fig. 5. Distribution of users at different time slots.

This dataset contains 17,621 trajectories with a total distance of about 1.2 million kilometers and a total duration of over 48,000 hours [42], [43]. Since this dataset recorded 182 users' outdoor trajectories in a broad range, we process it according to the features of users' activities. We first observed the activity tracks of 182 users and marked the longitude and latitude of the origin center coordinates [116.327544, 39.987317]. Then, we take this location as the central point and divide the area into three different scopes by setting the radius to $r = 1.0$, $r = 2.5$, and $r = 3.0$ kilometers. The division of the activities' scopes from different groups is shown in Fig. 4(a). Based on that, we construct three datasets with different characteristics by comprehensively considering the time, scopes and trajectories of users' activities. For each group of datasets, we select 40 users to construct our dataset \mathbf{U} and traverse their trajectories to find the ones within the areas under a scaled time series. We continuously collect users' data during 20 ($r = 1.0$ km), 60 ($r = 2.5$ km), and 120 ($r = 3.0$ km) consecutive time slots for each group, respectively. The trajectories of \mathbf{U} in different scopes are shown in Fig. 4(b) to (d). We can see that in the first group of users, the overlap of user activity trajectories is not obvious due to the small geographic location and time range. Then we expanded the activities' scopes of users while increasing the tracked time slots into 60 and 120. We found that the probability of trajectory overlap in the time slot increases. Specifically, we take the group in $r = 2.5$ during 60 consecutive time slots as an example to show the distribution of users in different slots in Fig. 5, which includes the initial locations in time slots 0, 20, 40, and 59 in Fig. 5(a), (b), (c), and (d). We found that the locations of users vary in different time slots,

however, the number of connected users will remain at a high level for edge servers with a high frequency of utility. Based on that, we simulate the edge computing network based on \mathbf{U} , and we set up 49 edge servers with the service range of 450 meters. We set the computing capacity of each server to range from 2 GHz to 5 GHz, and the data size of each service is 1 GB. The storage of each edge server ranges from 5 GB to 10 GB, which also denotes the number of services that can be placed on edge servers. Compared to the proposed online service provisioning strategy, three baselines are used.

- *USNP-only*: Services provisioning and updating without using the prediction information, and the decisions are only made by USNP.
- *USP-only*: Services provisioning and updating by using the prediction information, and the decisions are only made by USP.
- *O-OSP*: Online services provisioning and updating based on $O-OSP_\omega$ without considering ω steps prediction.

B. Experiment Results

1) *Average Total Delay Under Different Strategies*: We investigate the average total delays under these four strategies with four groups of users ($|\mathbf{U}| = 10$, $|\mathbf{U}| = 20$, $|\mathbf{U}| = 30$, and $|\mathbf{U}| = 40$) in different time scales. The results are shown in Figs. 7 to 8. Combined with the distribution characteristics of user activities in consecutive time slots for each group under different activity scopes 20 ($r = 1.0$), 60 ($r = 2.5$), and 120 ($r = 3$), respectively, we have the following observations.

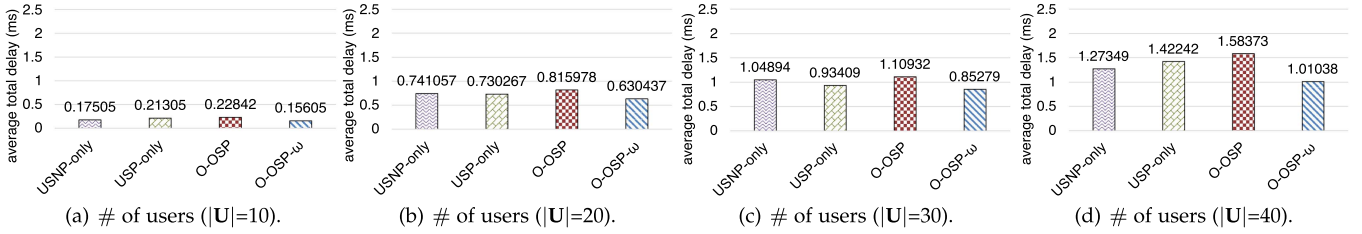


Fig. 6. Average total delay under different strategies of users with 20-step trajectory ($r = 1.0$ km).

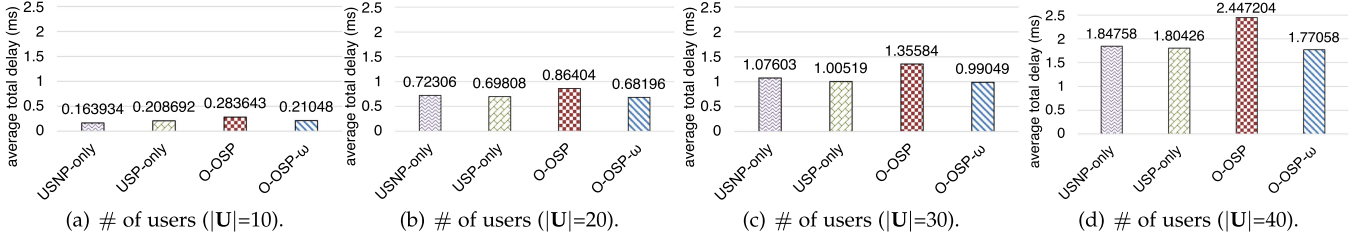


Fig. 7. Average total delay under different strategies of users with 60-step trajectory ($r = 2.5$ km).

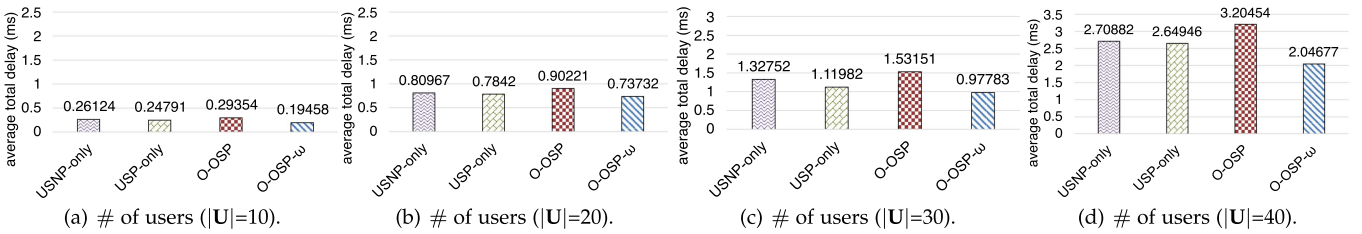
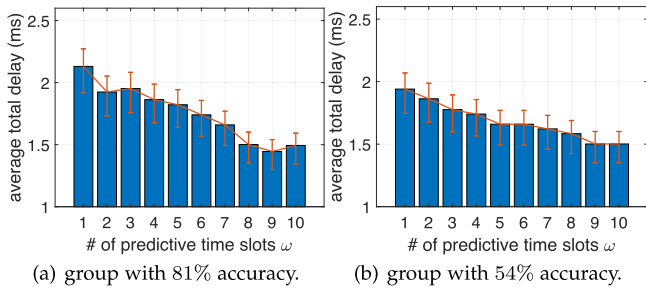
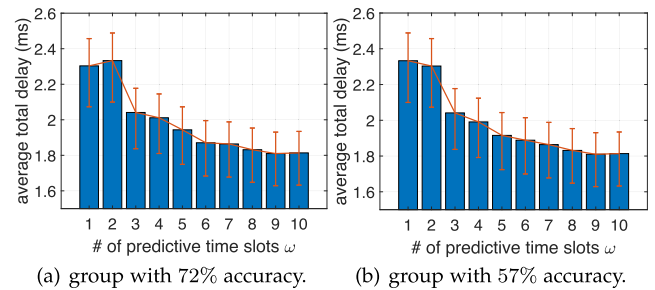


Fig. 8. Average total delay under different strategies of users with 120-step trajectory ($r = 3.0$ km).

i) For each group, the numbers of users in set \mathbf{U} affect the results of strategies. We analyze the average total delay for different groups of users under the same trajectory. Viewing the results in Figs. 6 to 8 as a whole, the tendencies of the average total delays in groups of 20-step, 60-step, and 120-step climb with the increasing number of users. The main reason for this situation is the resource competition problem caused by the increase of users in the same active area. For the group with the shortest trajectory steps (20-step), algorithm USNP-only has the largest fluctuation on the total delay, which reaches about 1.3 multiples on average. By comparison, the fluctuation of the average total delay under algorithms O-OSP and O-OSP- ω are relatively stable and slow. However, for different numbers of users ($|\mathbf{U}| = 10$, $|\mathbf{U}| = 20$, $|\mathbf{U}| = 30$, and $|\mathbf{U}| = 40$), the average total of algorithm O-OSP is the highest, which means that if you only focus on the direction of the users in the next step, there may be extra overhead by the prediction error or the repeated operations of users' trajectories. As shown in Fig. 6, O-OSP- ω can obtain the minimum delay for different groups of users, which means that considering the multi-step prediction can effectively help to avoid extra overhead being produced by the complex environment when users increase. The group with longest

trajectory steps (120-step) showed similar trends to the one with 20 steps. The group with 60 steps has a slight difference. As shown in Fig. 7(b), (c), and (d), O-OSP- ω has the lowest average total delays under the groups of 20, 30, and 40. Meanwhile, the average total delay decreases notably with an increasing number of users. However, as shown in Fig. 7(a), USNP-only obtains the lowest delay in the group with 10 users. On one hand, there are abundant resources when there are fewer users, which leads to deviation in the decision-making under the prediction trajectory. On the other hand, we found that the trajectories for the selected 10 users in group one hardly changed, which results in errors in the algorithms using the prediction information.

ii) For different groups, the activities' scopes and trajectories of users affect the results of strategies. Comparing Figs. 6 to 8, expanding the activities' scopes has no significant impact on the average total delay of users, which is reflected in that the tendency has not changed significantly on the whole. For each group, the total average latency of algorithm O-OSP is the largest. Compare to the other three strategies, the result of O-OSP- ω is the minimum. The delay of USNP-only is higher than that USP-only in different trajectory ranges when the numbers of users

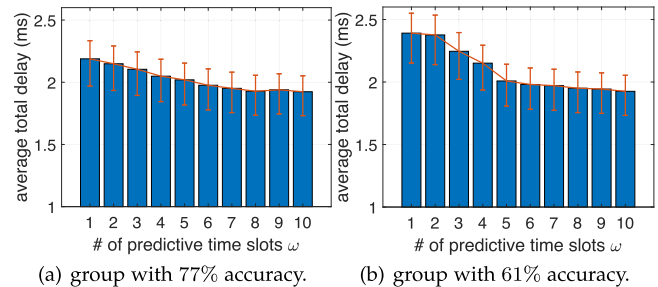

 Fig. 9. Average total delay of users with 20-step trajectory ($\omega \in [1, 10]$).

 Fig. 10. Average total delay of users with 60-step trajectory ($\omega \in [1, 10]$).

are 10, 20, and 30. The tendency is slightly different for the group with 40 users, which is presented in the results under these four strategies most obviously. Like the group with a small range of scale shown in Fig. 6(d), USP-only is better than USNP-only. This fluctuation is mainly due to the increased number of users moving in a small range, which will bring resource constraints. Once a service provisioning is improperly, this will result in a relatively strong impact. However, the average total delay using algorithm USP-only is indeed lower than USNP-only for users $|\mathcal{U}| = 40$ with a wide range of activities, 120-step trajectory. When the range of activities stays in the middle, the effects of these two algorithms (USNP-only and USP-only) are very close to the O-OSP $_{\omega}$, which means that they obtain better results in the current active range.

- iii) Prediction with ω slots in O-OSP $_{\omega}$ can effectively reduce the problem of service quality degradation caused by erratic activities of mobile users. As shown in Fig. 7(d), the average total delay of O-OSP becomes significantly higher than that of the other algorithms. In this case, besides the lowest average total delay of O-OSP $_{\omega}$, USNP-only and UPS-only can also achieve better performances. The reason is that the increase in delay is due to the scaling of users under limited resources. Especially in the case of the trajectories of users changing frequently, it may be inappropriate to determine the location of the service only by one step, which will affect the delay of other users. The simulation results show that our algorithm can reduce the average total delay of 28.7% (20-step trajectory, $r = 1.0$ km), 17.8% (60-step trajectory, $r = 2.5$ km), and 17.8% (120-step trajectory, $r = 3.0$ km) when comparing with baselines, respectively.

2) *Average Total Delay With Different ω Time Slots:* Based on the compared results above, we study the average total delay of O-OSP $_{\omega}$ with different predictive ω slots. We predict the trajectories of 40 users ($|\mathcal{U}|$) using the Social-LSTM model in multiple groups (20-step, 60-step, and 120-step), and we choose two sub-groups for each one with different accuracies. The results are shown in Figs. 9 to 11. Additionally, we have the following observations.

- i) The value of ω can influence the efficiency of O-OSP $_{\omega}$, and there are existing differences with the increase of ω in the tendency of the average total delay even if the prediction accuracies are similar. As shown in Figs. 9(a)


 Fig. 11. Average total delay of users with 120-step trajectory ($\omega \in [1, 10]$).

to 11(a), the accuracies of these groups are 81%, 77%, and 72% which means that the differences between any two groups' fluctuations within 10%. However, it is clear that the fluctuations between the average total delays under different prediction steps ω vary greatly. For the first group of 20-step, the average accuracy of prediction under different steps (ω) is the highest (81%). The initial change of ω , where increases from $\omega = 1$ to $\omega = 2$ occurred, resulted in a very significant drop in the average total delay. But when the steps increase to $\omega = 3$, the delay increased slightly. Then, there is an approximately linear decline from $\omega = 4$ to $\omega = 9$. When the number of the prediction step becomes too large where $\omega = 10$, the average total delay of users in the group with 20-step trajectory decreases. For the second group of 60-step, which are shown in Fig. 10(a) and 10(b), we have 72% and 57% percent accuracies for the comparative experiments. When the ω steps range from 1 to 9, the average total delay of users keeps decreasing. For each group, we can see that there is an obvious change between $\omega = 2$ and $\omega = 3$ which means that the initial change of ω has no effect on the delay, and the inflection point appears when $\omega = 3$. Then, there is an approximately linear decline from $\omega = 3$ to $\omega = 9$. When the slots scale into $\omega = 9$ and $\omega = 10$ ($\omega \geq 9$), the average total delay does not change notably. The reason for this is that the prediction of users' trajectories too far ahead of their movements may cause inaccurate results which may lead to invalid decisions. For group three with the widest range of activity scope (120-step), which are shown in Fig. 11(a) and 11(b), the average total delay decreases smoothly for the group with

77% accuracy. When the value of ω increases to 8, the subsequent increase of ω has little effect on the delay. From the analysis by comparing Fig. 9(a) to Fig. 11(a), we have that when the accuracy rates are higher and close, the limitation on the activity scope of users causes local anti-correlation between the number of prediction steps ω and the average total delay. Therefore, the total average delay under the O-OSP $_{\omega}$ strategy decreases in a range with the increasing value of ω , and the setting of ω is related to the characteristics of users and the prediction model. Comparing the three groups comprehensively, the number of predicted steps about $\omega = 9$ can obtain good results.

- ii) The accuracy of the chosen prediction model has little effect on the results of O-OSP $_{\omega}$. As shown in Figs. 9 to 11, we selected two sub-groups with large differences which can reach about 20% on average in the prediction accuracy under different scope of activities for analysis. When the number of prediction steps is extremely small, i.e., $\omega = 1$, the lower accuracy has little effect on the average total delay. In some groups, the delay of users decreases with accuracy. As per users with 60-step trajectories in Fig. 10, when the accuracy decreases to 57%, the average total delay under $\omega = 1$ is barely growing. For the users with 120-step trajectories in Fig. 11, compared the group with 61% accuracy to the 72% one, there is an obvious growth on the average total delay under $\omega = 1$ which is caused by the scaling activity scope. At the same time, the result of the decreased delay may also occur in the case of lower accuracy under a small prediction step, such as for users with 20-step trajectories in Fig. 9 when $\omega = 1$. When the number of prediction steps gradually change, the delays for these three groups decrease and the gap between sub-groups become narrower with the increase of ω . For the group with a wide range of activity which leads to higher uncertainty on the trajectories of users (120-step trajectory, $r = 3.0$ km), the total average delay with smaller prediction steps shows increasing tendencies, and the results are getting closer when the steps are over $\omega \geq 6$. For the other two groups of users (20-step trajectory, $r = 1.0$ km, and 120-step trajectory, $r = 3.0$ km) which are shown in Figs. 9 and 10, the average total delay under $\omega = 6$ in these two groups with 54% and 57% accuracies are also basically the same. Therefore, we have that even if the accuracy of the prediction model is imprecise, O-OSP $_{\omega}$ can still obtain a better result.

VII. CONCLUSION

In this article, we investigate the service provisioning and updating problem in a multi-user scenario by improving the performance of services within the long-term cost constraint. First, we decouple the original long-term optimization problem into a deterministic per-slot problem using Lyapunov optimization. Based on that, we propose two service updating decision strategies by considering the trajectory prediction conditions of users. Based on this, we design an online strategy by utilizing

the committed horizon control method while looking ahead to ω slots predictions. We prove the performance bound of our online strategy theoretically in terms of the trade-off between delay and cost. Finally, we conduct extensive experiments based on the Microsoft GPS trajectory dataset, and we demonstrate the superior performance of the proposed algorithm.

In our future work, we intend to investigate the applicability of our algorithm to edge service markets, focusing on the study of backup strategies and service reliability in dynamic environments. In addition, we plan to find feasible and effective solutions considering fairness.

REFERENCES

- [1] S. Tu, M. Waqas, S. U. Rehman, T. Mir, Z. Halim, and I. Ahmad, "Social phenomena and fog computing networks: A novel perspective for future networks," *IEEE Trans. Comput. Social Syst.*, vol. 9, no. 1, pp. 32–44, Feb. 2022.
- [2] M. Waqas, S. Tu, Z. Halim, S. U. Rehman, G. Abbas, and Z. H. Abbas, "The role of artificial intelligence and machine learning in wireless networks security: Principle, practice and challenges," *Artif. Intell. Rev.*, vol. 55, no. 7, pp. 5215–5261, 2022.
- [3] T. K. Dang, N. Mohan, L. Corneo, A. Zavodovski, J. Ott, and J. Kangasharju, "Cloudy with a chance of short RTTs: Analyzing cloud connectivity in the internet," in *Proc. 21st ACM Internet Meas. Conf.*, 2021, pp. 62–79.
- [4] Y. Chen, J. Wu, and B. Ji, "Virtual network function deployment in tree-structured networks," in *Proc. IEEE 26th Int. Conf. Netw. Protoc.*, 2018, pp. 132–142.
- [5] Y. Siriwardhana, P. Porambage, M. Liyanage, and M. Ylianttila, "A survey on mobile augmented reality with 5G mobile edge computing: Architectures, applications, and technical aspects," *IEEE Commun. Surv. Tut.*, vol. 23, no. 2, pp. 1160–1192, Secondquarter 2021.
- [6] F. A. Salaht, F. Desprez, and A. Lebre, "An overview of service placement problem in fog and edge computing," *ACM Comput. Surv.*, vol. 53, no. 3, pp. 1–35, 2020.
- [7] L. Wang, L. Jiao, T. He, J. Li, and H. Bal, "Service placement for collaborative edge applications," *IEEE/ACM Trans. Netw.*, vol. 29, no. 1, pp. 34–47, Feb. 2021.
- [8] K. Elgazzar, P. Martin, and H. S. Hassanein, "Empowering mobile service provisioning through cloud assistance," in *Proc. IEEE/ACM 6th Int. Conf. Utility Cloud Comput.*, 2013, pp. 9–16.
- [9] Y. Qiu, J. Liang, V. C. M. Leung, X. Wu, and X. Deng, "Online reliability-enhanced virtual network services provisioning in fault-prone mobile edge cloud," *IEEE Trans. Wireless Commun.*, vol. 21, no. 9, pp. 7299–7313, Sep. 2022.
- [10] J. Li, W. Liang, M. Huang, and X. Jia, "Reliability-aware network service provisioning in mobile edge-cloud networks," *IEEE Trans. Parallel Distrib. Syst.*, vol. 31, no. 7, pp. 1545–1558, Jul. 2020.
- [11] N. Yu, Q. Xie, Q. Wang, H. Du, H. Huang, and X. Jia, "Collaborative service placement for mobile edge computing applications," in *Proc. IEEE Glob. Commun. Conf.*, 2018, pp. 1–6.
- [12] Y. Mao, X. Shang, and Y. Yang, "Joint resource management and flow scheduling for SFC deployment in hybrid edge-and-cloud network," in *Proc. IEEE INFOCOM Conf. Comput. Commun.*, 2022, pp. 170–179.
- [13] L. Gu, Z. Chen, H. Xu, D. Zeng, B. Li, and H. Jin, "Layer-aware collaborative microservice deployment toward maximal edge throughput," in *Proc. IEEE INFOCOM Conf. Comput. Commun.*, 2022, pp. 71–79.
- [14] Z. Nezami, K. Zamanifar, K. Djemame, and E. Pournaras, "Decentralized edge-to-cloud load balancing: Service placement for the Internet of Things," *IEEE Access*, vol. 9, pp. 64983–65000, 2021.
- [15] G. Zhang, S. Zhang, W. Zhang, Z. Shen, and L. Wang, "Joint service caching, computation offloading and resource allocation in mobile edge computing systems," *IEEE Trans. Wireless Commun.*, vol. 20, no. 8, pp. 5288–5300, Aug. 2021.
- [16] H. Chen et al., "Mobility-aware offloading and resource allocation for distributed services collaboration," *IEEE Trans. Parallel Distrib. Syst.*, vol. 33, no. 10, pp. 2428–2443, Oct. 2022.
- [17] J. Xu, L. Chen, and P. Zhou, "Joint service caching and task offloading for mobile edge computing in dense networks," in *Proc. IEEE INFOCOM Conf. Comput. Commun.*, 2018, pp. 207–215.

- [18] Y. Ren, S. Shen, Y. Ju, X. Wang, W. Wang, and V. C. M. Leung, "EdgeMatrix: A. resources redefined edge-cloud system for prioritized services," in *Proc. IEEE INFOCOM Conf. Comput. Commun.*, 2022, pp. 610–619.
- [19] X. Shang, Y. Huang, Y. Mao, Z. Liu, and Y. Yang, "Enabling QoE support for interactive applications over mobile edge with high user mobility," in *Proc. IEEE INFOCOM Conf. Comput. Commun.*, 2022, pp. 1289–1298.
- [20] X. Wang, J. Ye, and J. C. S. Lui, "Decentralized task offloading in edge computing: A multi-user multi-armed bandit approach," in *Proc. IEEE INFOCOM Conf. Comput. Commun.*, 2022, pp. 1199–1208.
- [21] P. Han, Y. Liu, and L. Guo, "Interference-aware online multicomponent service placement in edge cloud networks and its AI application," *IEEE Internet Things J.*, vol. 8, no. 13, pp. 10557–10572, Jul. 2021.
- [22] R. Li, Z. Zhou, X. Chen, and Q. Ling, "Resource price-aware offloading for edge-cloud collaboration: A two-timescale online control approach," *IEEE Trans. Cloud Comput.*, vol. 10, no. 1, pp. 648–661, Jan.–Mar. 2022.
- [23] B. Liu, W. Zhang, W. Chen, H. Huang, and S. Guo, "Online computation offloading and traffic routing for UAV swarms in edge-cloud computing," *IEEE Trans. Veh. Technol.*, vol. 69, no. 8, pp. 8777–8791, Aug. 2020.
- [24] F. Liu, Z. Zhou, H. Jin, B. Li, B. Li, and H. Jiang, "On arbitrating the power-performance tradeoff in SaaS clouds," *IEEE Trans. Parallel Distrib. Syst.*, vol. 25, no. 10, pp. 2648–2658, Oct. 2014.
- [25] W. Fang, X. Yao, X. Zhao, J. Yin, and N. Xiong, "A stochastic control approach to maximize profit on service provisioning for mobile cloudlet platforms," *IEEE Trans. Syst. Man Cybern. Syst.*, vol. 48, no. 4, pp. 522–534, Apr. 2018.
- [26] Y. Qi, L. Pan, and S. Liu, "A Lyapunov optimization-based online scheduling algorithm for service provisioning in cloud computing," *Future Gener. Comput. Syst.*, vol. 134, pp. 40–52, 2022.
- [27] Z. Ning et al., "Distributed and dynamic service placement in pervasive edge computing networks," *IEEE Trans. Parallel Distrib. Syst.*, vol. 32, no. 6, pp. 1277–1292, Jun. 2021.
- [28] T. Kim et al., "MoDEMS: Optimizing edge computing migrations for user mobility," in *Proc. IEEE INFOCOM Conf. Comput. Commun.*, 2022, pp. 1159–1168.
- [29] Y. Zeng, Y. Huang, Z. Liu, and Y. Yang, "Online distributed edge caching for mobile data offloading in 5G networks," in *Proc. IEEE/ACM 28th Int. Symp. Qual. Serv.*, 2020, pp. 1–10.
- [30] Z. Li, C. Jiang, and J. Lu, "Distributed service migration in satellite mobile edge computing," in *Proc. IEEE Glob. Commun. Conf.*, 2021, pp. 1–6.
- [31] E. Liu, X. Deng, Z. Cao, and H. Zhang, "Design and evaluation of a prediction-based dynamic edge computing system," in *Proc. IEEE Glob. Commun. Conf.*, 2018, pp. 1–6.
- [32] Y. Jin, L. Jiao, Z. Qian, S. Zhang, and S. Lu, "Learning for learning: Predictive online control of federated learning with edge provisioning," in *Proc. IEEE INFOCOM Conf. Comput. Commun.*, 2021, pp. 1–10.
- [33] H. Ma, Z. Zhou, and X. Chen, "Leveraging the power of prediction: Predictive service placement for latency-sensitive mobile edge computing," *IEEE Trans. Wireless Commun.*, vol. 19, no. 10, pp. 6454–6468, Oct. 2021.
- [34] T. Ouyang, Z. Zhou, and X. Chen, "Follow me at the edge: Mobility-aware dynamic service placement for mobile edge computing," *IEEE J. Sel. Areas Commun.*, vol. 36, no. 10, pp. 2333–2345, Oct. 2018.
- [35] L. Ale, N. Zhang, X. Fang, X. Chen, S. Wu, and L. Li, "Delay-aware and energy-efficient computation offloading in mobile-edge computing using deep reinforcement learning," *IEEE Trans. Cogn. Commun. Netw.*, vol. 7, no. 3, pp. 881–892, Sep. 2021.
- [36] Z. Xiao et al., "Multi-objective parallel task offloading and content caching in D2D-aided MEC networks," *IEEE Trans. Mobile Comput.*, vol. 22, no. 11, pp. 6599–6615, Nov. 2023.
- [37] S. Lu, J. Wu, J. Shi, P. Lu, J. Fang, and H. Liu, "A dynamic service placement based on deep reinforcement learning in mobile edge computing," *Network*, vol. 2, no. 1, pp. 106–122, 2022.
- [38] T. Taleb, A. Ksentini, and P. A. Frangoudis, "Follow-me cloud: When cloud services follow mobile users," *IEEE Trans. Cloud Comput.*, vol. 7, no. 2, pp. 369–382, Apr.–Jun. 2019.
- [39] B. Gao, Z. Zhou, F. Liu, and F. Xu, "Winning at the starting line: Joint network selection and service placement for mobile edge computing," in *Proc. IEEE INFOCOM Conf. Comput. Commun.*, 2019, pp. 1459–1467.
- [40] M. J. Neely, "Stochastic network optimization with application to communication and queueing systems," *Synth. Lectures Commun. Netw.*, vol. 3, no. 1, pp. 1–211, 2010.
- [41] J. Comden, S. Yao, N. Chen, H. Xing, and Z. Liu, "Online optimization in cloud resource provisioning: Predictions, regrets, and algorithms," in *Proc. ACM Meas. Anal. Comput. Syst.*, vol. 3, no. 1, pp. 1–30, 2019.
- [42] Y. Zheng, Q. Li, Y. Chen, X. Xie, and W. Y. Ma, "Understanding mobility based on GPS data," in *Proc. 10th Int. Conf. Ubiquitous Comput.*, 2008, pp. 312–321.
- [43] Y. Zheng, L. Zhang, X. Xie, and W. Y. Ma, "Mining interesting locations and travel sequences from GPS trajectories," in *Proc. 18th Int. Conf. World Wide Web*, 2009, pp. 791–800.
- [44] G. Qiao, S. Leng, S. Maharjan, Y. Zhang, and N. Ansari, "Deep reinforcement learning for cooperative content caching in vehicular edge computing and networks," *IEEE Internet Things J.*, vol. 7, no. 1, pp. 247–257, Jan. 2020.
- [45] S. Wang, Y. Guo, N. Zhang, P. Yang, A. Zhou, and X. Shen, "Delay-aware microservice coordination in mobile edge computing: A reinforcement learning approach," *IEEE Trans. Mobile Comput.*, vol. 20, no. 3, pp. 939–951, Mar. 2021.



Shuaibing Lu (Member, IEEE) received the PhD degree in computer science and technology from Jilin University, Changchun, in 2019. She is currently a lecturer with the Faculty of Information Technology, Beijing University of Technology. She is supported by the China Scholarship Council as a visiting scholar supervised by Prof. From 2016 to 2018, she was with the Department of Computer and Information Sciences, Temple University. Her research interests include distributed computing, cloud computing, and edge computing.



Jie Wu (Fellow, IEEE) is currently the director of the Center for Networked Computing and Laura H. Carnell professor with Temple University. He is the director of international affairs with the College of Science and Technology. He was the Chair of Department of Computer and Information Sciences from the summer of 2009 to summer of 2016 and associate vice provost for international affairs from the fall of 2015 to summer of 2017. Prior to joining Temple University, he was a program director with the National Science Foundation and was a distinguished professor with Florida Atlantic University. He regularly publishes in scholarly journals, conference proceedings, and books. His research interests include mobile computing and wireless networks, routing protocols, cloud and green computing, network trust and security, and social network applications. He serves on several editorial boards, including *IEEE Transactions on Service Computing* and *Journal of Parallel and Distributed Computing*. Dr. Wu was the general co-chair for IEEE MASS 2006, IEEE IPDPS 2008, IEEE ICDCS 2013, ACM MobiHoc 2014, ICPP 2016, and IEEE CNS 2016, and program co-chair for IEEE INFOCOM 2011 and CCF CNCC 2013. He was an IEEE Computer Society distinguished visitor, ACM distinguished speaker, and chair for the IEEE Technical Committee on Distributed Processing. Dr. Wu is a CCF distinguished speaker. He was the recipient of the 2011 China Computer Federation (CCF) Overseas Outstanding Achievement Award.



Pengfan Lu received the BSc degree in computer science and technology with the Harbin University of Science and Technology. He is currently working toward the MSc degree in computer science with the Faculty of Information Technology, Beijing University of Technology. His research interests include cloud computing and edge computing.



Ning Wang received the BE degree from the School of Physical Electronics, University of Electronic Science and Technology of China, Chengdu, Sichuan, China, in 2013, and the PhD degree from the Department of Computer and Information Sciences, Temple University, Philadelphia, PA, USA, in 2018. He is currently an Assistant Professor with the Department of Computer Science, Rowan University, Glassboro, NJ. He is focuses on communication and computation optimization problems in Internet-of-Things systems, and operation optimization in smart cities applica-

tions. He has authored or coauthored nearly 30 papers in high-impact networking conferences and journals, such as, IEEE International Conference on Distributed Computing Systems, IEEE INFOCOM, IEEE/ACM International Symposium on Quality of Service, *IEEE Transactions on Big Data*, and *Journal of Parallel and Distributed Computing*. He was a program committee member for top international conferences such as IEEE International Conference on Distributed Computing Systems and IEEE Wireless Communications and Networking Conference, and reviewers for premier journals, such as *IEEE Transactions on Parallel and Distributed Systems*, *IEEE Transactions on Wireless Communications*, *IEEE Transactions on Mobile Computing*, *IEEE Transactions on Intelligent Transportation Systems*, *ACM Transactions on Internet Technology*, *IEEE Transactions on Intelligent Transportation Systems*, and *IEEE Transactions on Services Computing*.



Haiming Liu received the BS degree in computer science and technology, the MS degree in computer software and theory, and the PhD degree in computer science and technology (bioinformatics) from Jilin University, Changchun, in 2012, 2015, and 2019, respectively. He is currently a lecturer with the School of Software Engineering, Beijing Jiaotong University. He is a member of Chinese Association for Artificial Intelligence. His research interests include edge computing, data mining, and bioinformatics.



Juan Fang (Member, IEEE) received the MS degree from the Jilin University of Technology, Changchun, China, in 1997, and the PhD degree from the College of Computer Science, Beijing University of Technology, Beijing, China, in 2005. In 1997, she joined the College of Computer Science, Beijing University of Technology. From 2015, she has been a professor with the Beijing University of Technology. Her research interests include high performance computing, edge computing, and Big Data technology.