

Robust Resource Provisioning in Time-Varying Edge Networks

Ruozhou Yu
North Carolina State University
Raleigh, NC

Guoliang Xue, Yinxin Wan
Arizona State University
Tempe, AZ

Jian Tang
Syracuse University
Syracuse, NY

Dejun Yang
Colorado School of Mines
Golden, CO

Yusheng Ji
National Institute of Informatics
Tokyo, Japan

ABSTRACT

Edge computing is one of the revolutionary technologies that enable high-performance and low-latency modern applications, such as smart cities, connected vehicles, etc. Yet its adoption has been limited by factors including high cost of edge resources, heterogeneous and fluctuating demands, and lack of reliability. In this paper, we study resource provisioning in edge computing, taking into account these different factors. First, based on observations from real demand traces, we propose a time-varying stochastic model to capture the time-dependent and uncertain demand and network dynamics in an edge network. We then apply a novel robustness model that accounts for both expected and worst-case performance of a service. Based on these models, we formulate edge provisioning as a multi-stage stochastic optimization problem. The problem is NP-hard even in the deterministic case. Leveraging the multi-stage structure, we apply nested Benders decomposition to solve the problem. We also describe several efficiency enhancement techniques, including a novel technique for quickly solving the large number of decomposed subproblems. Finally, we present results from real dataset-based simulations, which demonstrate the advantages of the proposed models, algorithm and techniques.

CCS CONCEPTS

• **Networks** → **Network resources allocation; Traffic engineering algorithms; Network management; Mobile networks.**

KEYWORDS

Edge Computing, Resource Allocation, Time-Varying, Robustness, Multi-Stage Stochastic Optimization

ACM Reference Format:

Ruozhou Yu, Guoliang Xue, Yinxin Wan, Jian Tang, Dejun Yang, and Yusheng Ji. 2020. Robust Resource Provisioning in Time-Varying Edge Networks. In *The Twenty-first ACM International Symposium on Theory, Algorithmic Foundations, and Protocol Design for Mobile Networks and Mobile Computing (MobiHoc '20)*, October 11–14, 2020, Boston, MA, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3397166.3409146>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions.acm.org.
MobiHoc '20, October 11–14, 2020, Boston, MA, USA

© 2020 Association for Computing Machinery.
ACM ISBN 978-1-4503-8015-7/20/10...\$15.00
<https://doi.org/10.1145/3397166.3409146>

1 INTRODUCTION

Advances in machine learning, computer vision, big data and more have enabled many new applications that employ connected intelligence to improve human lives, some of the most promising ones including smart home, smart city, and autonomous driving, to name a few. These applications gather the massive data generated by networked devices, and use centralized computing for data analytics and decision making. With more and more devices connected via the Internet-of-Things, however, our existing communication and computing infrastructures are being challenged. How to meet application performance goals while facing excessive load is becoming a critical problem both technically and commercially, which hinders the adoption of new life-changing applications.

Our current Internet-based cloud platforms fail to provide high-throughput, low-latency and reliable services due to unpredictability of the Internet. This has motivated the development of edge computing, where computing resources are provided via distributed nodes in the edge network. Edge computing enables services deployment proximal to end devices and users, both improving throughput and reducing latency. Yet edge computing has its own limitations, such as high cost, limited capacity, unreliable networks, etc. Efficient utilization of edge resources thus becomes a crucial problem.

In this paper, we study resource provisioning for edge services. We consider an integrated provisioning problem that involves service deployment, network provisioning, and simultaneously estimating and optimizing service performance. Different factors are accounted for during provisioning, including cost budget of the service provider, limited network resources, service robustness, and fluctuating demands and network condition. For the last consideration, we base our study on analysis of real-world datasets. We observe that some applications have very heterogeneous demands across both time and space. Moreover, in many cases the demands follow temporal-spatial repeating patterns, which is commonly due to the repeating behaviors of its users and devices. Such patterns can be crucial in performance enhancement and resource utilization. We propose a time-varying stochastic model to account for these patterns, which allows finer-grained resource provisioning. We further consider service performance robustness in terms of average user quality-of-service (QoS). We apply a balanced robustness metric, named Conditional Value-at-Risk (CVaR), which can flexibly trade-off between the expected and worst-case service performance.

Based on these models, we formulate robust edge provisioning as a three-stage stochastic optimization problem. The problem is a large-scale mixed integer linear program (MILP), and is NP-hard even in the deterministic case. We propose an optimal algorithm

based on *nested Benders decomposition*, which utilizes the decomposable structure of the problem. The original Benders decomposition suffers from slow convergence. We then describe efficiency enhancement techniques, including both established ones that apply to our problem, and a novel analytical solution to the decomposed subproblems. To validate our study, we have conducted simulations based on datasets synthesized from real data traces. The results show that our model enables fine-grained and robust provisioning of edge resources, the proposed algorithm generates superior solutions over heuristic approaches, and the described efficiency enhancement techniques are effective in reducing convergence time. To summarize, our contributions are as follows:

- We observe temporal-spatial repeating demand patterns from real-world traces, and formulate a three-stage edge resource provisioning problem with cost, resource, and robustness considerations. Edge provisioning with time-varying and temporal-spatial repeating demands has yet been studied to the best of our knowledge.
- We propose a nested Benders decomposition-based optimal approach to this NP-hard problem, as well as several efficiency enhancement techniques.
- We present results from real data trace-based simulations, which show the advantages of the proposed models, algorithm, and techniques.

The rest of this paper is organized as follows. Sec. 2 presents our observations and our system model. Sec. 3 presents our problem and states its complexity. Sec. 4 presents our optimal algorithm and the efficiency enhancement techniques. Sec. 5 presents our simulation results. Sec. 6 discusses related work in the literature. Finally, Sec. 7 concludes this paper.

2 SYSTEM MODEL AND OVERVIEW

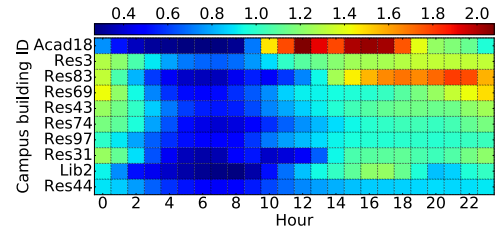
2.1 Edge System Model

The edge network is modeled as a directed graph $G = (N, L)$, where N is the set of nodes, and L is the set of links between nodes. The network consists of *computing nodes* (hosts) denoted by $H \subseteq N$, and *access points* (APs) denoted by $A \subseteq N$. Each link $l \in L$ has a bandwidth capacity b_l . We assume the network is strongly connected. A weakly connected or disconnected network can be made strongly connected by adding artificial zero-bandwidth links.

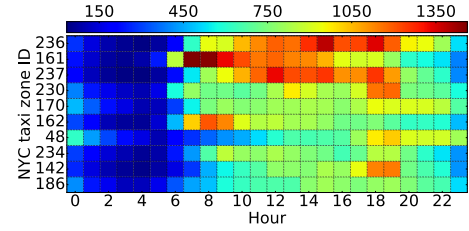
Three parties exist in the system. The *network manager* (NM) manages the network, including all network nodes and links. The *edge computing manager* (ECM) manages the computing resources on computing nodes. Finally, the *edge service provider* (ESP) operates an *edge service*. The ESP's goal is to deploy its service in the edge network to serve external user demands. We consider an edge service receiving and processing continuous data streams from APs. Typical examples include crowdsensing services that aggregate and analyze sensing data from mobile users, smart transportation services that receive data from sensors and vehicles and perform traffic scheduling, and emergency services that monitor and extract eventual information. The service can be deployed on one or multiple distributed computing nodes. However, as edge resources are expensive, universal service deployment may be impractical. Formally, let us use c_h to denote the cost of deploying and operating the edge service on node $h \in H$. The cost budget of the ESP is C .

2.2 Edge Demand Model

While demand patterns differ across applications, a common characteristic of edge demand is its frequent temporal-spatial dynamics. Such dynamics usually follow time-varying and repeating patterns over time. This could be due to many factors such as repeating user activities, user mobility, and duty-cycling services, and it can be observed on both small-scale and large-scale systems. For example, Fig. 1(a) shows the average-per-hour number of mobile devices connected to top-10 popular building APs in a campus network from 09/2002 to 09/2003 [15]. Fig. 1(b) shows the average number of passengers dropped-off by Yellow Taxi at top-10 popular taxi zones in NYC in 2018 [27]. Though measuring different data, both datasets show the different demands in different hours over a day, which most likely attributes to the repeating daily activities of users.



(a) Top 10 Dartmouth College buildings: average connected devices



(b) Top 10 New York City (NYC) taxi zones: average passenger drop-offs

Figure 1: Time-varying demands from two different datasets: (a) average number of devices connected to campus APs, and (b) average passenger drop-offs at NYC taxi zones.

Motivated by the observation, we adopt a time-varying demand model. We divide time into T time slots, each representing a demand pattern observed over time. For simplicity, we assume all slots are of equal lengths and are successively repeating, and denote the course of time over all slots as a *period*. Our model and solution can be trivially extended to handle time slots with different lengths and repeating patterns. From Fig. 1, demand distributions can be neither identical nor independent across APs in one time slot, or across different slots at the same AP. Specifically, demand at AP $a \in A$ in slot $t = 1 \dots T$ is assumed to follow a (possibly unknown) distribution, which we denote by random variable $\delta_{t,a}$; the total demand in t is defined as $\delta_t = \sum_a \delta_{t,a}$. This reflects the fact that even in the same slot, the exact demand may still fluctuate across periods due to unattended factors such as special events, seasonal changes, etc. Under this model, provisioning is done for each slot across different periods, instead of being fixed at all time.

2.3 Network Provisioning Model

Network has a critical impact on edge service performance. Ideally, the ESP wishes to attain optimal network performance via frequent

or real-time network optimization. However, the network is managed by the NM (e.g., an ISP) and hence out of ESP's control. We assume the NM can process network optimization requests from the ESP, but not frequently or in real-time due to the complexity of network management and the large number of supported services.

Based on the time-varying demand model, the ESP employs the following strategy. During initial provisioning, the ESP submits a network optimization request for each time slot with demand information. The NM provisions a dedicated networking policy for each slot based on the requests. After the service starts operating, the NM updates its policy based on the time slots. If the demand patterns change, the ESP can re-submit those requests, and the NM will derive new policies and operate based on them thereafter.

Routing. The NM's network policy can greatly impact the service performance. Ideally, flow-based optimization (e.g., min-cost multi-commodity flow [1]) can lead to optimal routing and bandwidth allocation for maximum service performance. Practically, the high complexity of flow-based algorithms discourages their usage in real networks. Fortunately, recent advances in *pathbook routing* have shown that a small set of carefully selected paths (e.g., 3–5 paths per commodity) can achieve close-to-optimal throughput in real networks, while traffic engineering overhead can be drastically reduced compared to flow-based algorithms [16]. Since the NM may serve many ESP requests at the same time, we believe that pathbook routing is suitable as an efficient yet adequate routing solution for the NM. Formally, we use $P = \bigcup \{P_{a,h} \mid a \in A, h \in H\}$ to denote the pathbook used by the NM. We only consider paths that connect APs to hosts for a specific ESP. Based on ESP requests, the NM derives bandwidth allocation across paths for each time slot, subject to the available bandwidth on network links.

Performance of an edge service is largely dependent on network performance of data streams. We assume ESP's goal is to optimize the QoS of its service, notably the network delay, when serving external demands. Network conditions may fluctuate due to congestion, failures and maintenance, malicious traffic, etc. To model such dynamics, we use random variable $d_{t,l}$ to denote the *delay* of a link $l \in L$ in time slot t . We note that link delay is independent of demand due to explicit demand allocation; see next section.

3 PROBLEM STATEMENT AND COMPLEXITY

The initial edge provisioning consists of the following steps. First, the ESP selects a set of computing nodes to deploy the service. Since service deployment incurs manual cost for setup and testing, it is *fixed* once the provisioning is done; live service migration is considered harmful due to interrupted service to users, and hence is not considered in this paper. Second, the ESP submits requests to the NM for network optimization based on demand and network statistics. After the provisioning is done, the ESP is presented with a service deployment plan, as well as a set of paths with pre-determined per-time slot bandwidth quotas. Unlike service deployment by the ECM, the NM can change bandwidth quotas by simply adjusting weights in the network without interrupting the service. Hence the bandwidth quotas can be updated in each time slot. In operations, the ESP adjusts demand allocation across paths in real-time, subject to the bandwidth quotas of the paths, to continuously serve user demands with optimized network performance.

3.1 Basic Problem Formulation

3.1.1 Service Deployment (SD). The first stage, service deployment, is formally denoted by variables $X = \{x(h) \in \{0, 1\} \mid h \in H\}$, where $x(h) = 1$ means that edge node h hosts a service instance and vice versa. The deployment is bounded by the cost budget of the ESP:

$$\sum_{h \in H} c_h x(h) \leq C. \quad (1)$$

Further, the service must have at least one instance deployed:

$$\sum_{h \in H} x(h) \geq 1. \quad (2)$$

3.1.2 Network Provisioning (NPR). Given fixed SD, the next step is to provision the network for each time slot. Define $Y = \bigcup_{t=1}^T Y_t$, where $Y_t = \{y(t, p) \geq 0 \mid p \in P\}$ is the set of variables for *pathbook bandwidth allocation* for slot t . First, the allocation is constrained by the deployed service instance, i.e., a path can have positive bandwidth only if it connects an access point to a deployed instance:

$$y(t, p) \leq b_p^{\max} x(h_p), \quad \forall t, p \in P, \quad (3)$$

where $b_p^{\max} = \min_{l \in p} \{b_l\}$ is the bottleneck bandwidth of the path, and h_p is the host destination of path p . Further, in each t , all paths using a link must obey the link's bandwidth capacity:

$$\sum_{p \in P: l \in p} y(t, p) \leq b_l, \quad \forall t, l \in L. \quad (4)$$

3.1.3 Network Estimation (NE). In normal operations, the ESP would flexibly allocate demands across paths subject to the pre-determined bandwidth allocation, in order to optimize instantaneous network performance. This information can also be used in initial provisioning, when the ESP estimates the *expected network performance* for a bandwidth allocation plan, and optimizes both service deployment and bandwidth allocation accordingly. Formally, define $Z = \bigcup_t Z_t$, where $Z_t = \{z(t, p) \geq 0 \mid p \in P\}$ is the set of (random) variables denoting demand allocation on path p in time t . The allocation is upper bounded by the bandwidth on each path:

$$z(t, p) \leq y(t, p), \quad \forall t, p \in P. \quad (5)$$

The allocation must also serve all demand from each AP:

$$\sum_{p \in P_a} z(t, p) \geq \delta_{t,a}, \quad \forall t, a \in A, \quad (6)$$

where $P_a = \bigcup_{h \in H} P_{a,h}$ contains all paths from a to a host.

3.1.4 Network Objective. With Constraint (6) serving all user demands, the ESP's goal is to optimize network QoS, i.e., transmission delay of all user demands. Given the stochastic model of demands and delays, we use D_t to denote the weighted average delay of all demands at time t , which is also a random variable, where

$$D_t \triangleq \frac{1}{\delta_t} \sum_{p \in P} d_{t,p} z(t, p). \quad (7)$$

and $d_{t,p} = \sum_{l \in p} d_{t,l}$ is the delay of a path. Consistent service performance is preferred, hence ESP can seek to optimize the highest delay across time. Informally, the problem can be expressed as

$$\min_{X \in \mathcal{F}} \max_t \{D_t\}, \quad (8)$$

where $X = (X, Y, Z)$, and \mathcal{F} is the feasibility region by (1)–(6).

3.2 Stochastic Optimization and Robustness

Given the stochastic nature of D_t , a common practice is to optimize the *expectation* of the objective. However, expectation optimization has limited ability in handling infrequent adverse scenarios [36]. For the provider, an elastic and fault-tolerant system produces more consistent revenue and lowers business risk in the long run. For

the users, a service with more stable performance leads to better user experience and higher satisfaction. Both perspectives coincide on the importance of *robustness* in performance optimization.

To counter performance fluctuation, we use a model called the *Conditional Value-at-Risk* (CVaR), a common risk measure in financial risk management. CVaR is widely used to measure the expected investment loss that an investor faces when market undergoes unfavorable variations. Formally, denote random variable \mathbf{R} as the investor's loss, then given a confidence value $\alpha \in [0, 1]$, define

$$\text{VaR}_\alpha(\mathbf{R}) \triangleq \min\{r \mid \Pr[\mathbf{R} \leq r] \geq \alpha\}, \quad (9)$$

$$\text{CVaR}_\alpha(\mathbf{R}) \triangleq \mathbb{E}[\mathbf{R} \mid \mathbf{R} \geq \text{VaR}_\alpha(\mathbf{R})]. \quad (10)$$

$\mathbb{E}(\cdot)$ denotes the expectation of a random variable. Eq. (9) defines the *Value-at-Risk* (VaR) with confidence α , representing the minimum value such that with α confidence, the investment loss will not exceed this value. Eq. (10) then defines the CVaR as the expected value of loss in all scenarios where the loss exceeds the VaR, i.e., the expected loss in the worst $(1 - \alpha)$ scenarios.

Notably, CVaR defines a trade-off between two objectives, controlled by α . If $\alpha = 0$, Eq. (10) equals to expectation of \mathbf{R} over the entire distribution. If $\alpha \rightarrow 1$, CVaR converges to the maximum loss over all possible scenarios, and minimizing CVaR equals to *minimax* in robust optimization. While expectation lacks robustness at all, minimax addresses only the worst case; in practice, minimax could reduce worst-case loss by an arbitrarily small amount, but drastically increase loss in all other cases. CVaR avoids both issues.

3.3 CVaR-based Formulation

Back to our problem, the ESP may seek to optimize both the expected average delay and CVaR of the average delay. To strike a balance, we formulate a *scalarized* two-objective problem:

$$\min_{X \in \mathcal{F}} \max_t \left\{ \rho_1 \cdot \mathbb{E}[D_t] + \rho_2 \cdot \text{CVaR}_\alpha(D_t) \right\}, \quad (11)$$

Parameters ρ_1 and ρ_2 define the trade-off between expectation and robustness, and can be tuned based on the ESP's actual objective.

In Eq. (10), the computation of CVaR requires first computing VaR, which itself does not have a closed form. Fortunately, Rockafellar and Uryasev have shown in [24] that the CVaR can be directly computed using the following formula:

$$\text{CVaR}_\alpha(D_t) = \min_{r(t) \in \mathbb{R}} \left\{ r(t) + \frac{1}{1-\alpha} \mathbb{E}[(D_t - r(t))^+] \right\}, \quad (12)$$

where $(D_t - r(t))^+ = \max\{D_t - r(t), 0\}$, and \mathbb{R} is the real number set. Note that in (12), the VaR is automatically computed as the minimizing value of $r(t)$. Merging (12) into Program (11), we have

$$\min_{X \in \mathcal{F}, R} \max_t \left(\rho_1 \mathbb{E}[D_t] + \rho_2 \left(r(t) + \frac{1}{1-\alpha} \mathbb{E}[(D_t - r(t))^+] \right) \right). \quad (13)$$

where $R = \{r(t)\}$ is the set of VaR variables for all time slots.

3.4 Sample Average Approximation (SAA)

Program (13) is hard to solve even with known distributions, since the solution cannot be derived in a closed form. Hence for each t , we use a set of independently and identically distributed (i.i.d.) samples $\mathcal{S}_t = (S_t^1, \dots, S_t^K)$ to approximate the distributions, where each $S_t^k = \{\delta_{t,a}^k \mid a \in A\} \cup \{d_{t,l}^k \mid l \in L\}$ contains a realization of the demands and delays in time t . With fixed SD and NPR, the optimal average delay of a sample can be computed by solving (5)–(8).

Let $Z_t^k = \{z(t, k, p)\}$ be the demand allocation variables for sample k of time slot t , and D_t^k be the optimal average delay (under (5)–(6)) for k . We use $\mathcal{Z}_t = \bigcup_k Z_t^k$ and $\mathcal{Z} = \bigcup_t \mathcal{Z}_t$ to denote the extended rate allocation sets over all samples of a time slot and all slots, respectively. The expected average delay $\mathbb{E}[D_t]$ is approximated by the sample average $\frac{1}{K} \sum_{k=1}^K D_t^k$. Similarly, $\text{CVaR}_\alpha(D_t)$ is approximated by $\min_{r(t) \in \mathbb{R}} \left\{ r(t) + \frac{1}{1-\alpha} \frac{1}{K} \sum_{k=1}^K (D_t^k - r(t))^+ \right\}$. To completely re-formulate Program (13) using sample averages, we define auxiliary variables $W = \bigcup_t W_t$, where $W_t = \{w(t, k) \geq 0 \mid k = 1 \dots K\}$. We also linearize the maximum term in the objective function (13) using an additional variable $D \in \mathbb{R}$. Then Program (13) can be re-formulated as a mixed integer linear program (MILP):

$$\min_{\substack{X, Y, \mathcal{Z}, \\ R, W, D}} D \quad (14a)$$

$$\begin{aligned} D &\geq \rho_1 \frac{1}{K} \sum_{k=1}^K D_t^k + \\ \text{s.t.} \quad &\rho_2 \left(r(t) + \frac{1}{1-\alpha} \frac{1}{K} \sum_{k=1}^K w(t, k) \right), \forall t; \end{aligned} \quad (14b)$$

$$w(t, k) \geq D_t^k - r(t), \quad \forall t, k; \quad (14c)$$

$$(1)-(4), \text{ and } \forall t, k, (5)-(6).$$

We call this the **Robust Edge Provisioning (REP)** problem. It is well-known that the SAA problem converges arbitrarily close to the original problem with a sufficient number of i.i.d. samples [14].

3.5 Computational Complexity

THEOREM 3.1. *REP is NP-hard.* \square

PROOF. Neglecting stochastics, the problem can be reduced from Knapsack, where \mathbf{n} items each with a cost c_i and a value v_i are to be put into a knapsack with maximum cost C and minimum value V . Consider an instance of REP with \mathbf{n} edge computing nodes in H , all connected to a single AP $a \in A$, each with deployment cost c_i and bandwidth v_i . The ESP has a budget of C , and a has a constant demand of V . Then deciding if the REP instance is feasible is equal to solving the Knapsack instance. REP's NP-hardness follows. \square

4 REP WITH MULTI-STAGE OPTIMIZATION

Due to problem NP-hardness and its scale, a decomposition-based algorithm is proposed based on the structure of the problem.

4.1 Problem Decomposition

From its definition, the problem can be vertically decomposed into three stages. At the later stages (NPR and NE), the corresponding subproblems can further be horizontally decomposed. The first-stage SD problem is an MILP:

$$O \triangleq \min_{X, D} D \text{ s.t. } \forall t, D \geq \mathcal{P}_t(X); (1)-(2). \quad (15)$$

The second-stage NPR subproblem given fixed X is an LP:

$$\begin{aligned} \mathcal{P}_t(X) &\triangleq \min_{Y_t, r(t)} \rho_2 r(t) + \frac{1}{K} \sum_k Q_{t,k}(Y, R) \\ \text{s.t.} \quad &(3)-(4). \end{aligned} \quad (16)$$

The third-stage NE subproblem given (Y, R) is also an LP:

$$\begin{aligned} Q_{t,k}(Y, R) &\triangleq \min_{Z, W} \rho_1 \cdot D_t^k + \rho_2 \cdot \frac{1}{1-\alpha} \cdot w(t, k) \\ \text{s.t.} \quad &(14c); \forall t, k, (5)-(6). \end{aligned} \quad (17)$$

We use $Q_{t,k}(Y, R)$ to omit index t at variables Y and R as function arguments. Similar notations are used hereafter to simplify notations. Eqs. (15)–(17) show that REP has a natural vertical-horizontal three-stage decomposition, where each $\mathcal{P}_t(X)$ subproblem is independent from its siblings, and so is each $Q_{t,k}(Y, R)$. Note that variables R and W , though both used to formulate CVaR, have been separated by the NPR–NE decomposition to enable full independence among $Q_{t,k}(Y, R)$ subproblems. This decomposition motivates us to adopt a decomposition-based algorithm to address the REP problem.

4.2 Benders Decomposition Basics

Benders decomposition is an iterative algorithm for solving multi-stage decomposable optimization problems [3]. Consider a general bounded two-stage problem in the following equal forms (assume that for any $X \in \mathcal{X}$ the problem is bounded):

$$\min_{X \in \mathcal{X}, Y \geq 0} \{dX + eY \mid AX + BY \geq f\} \quad (18)$$

$$= \min_{X \in \mathcal{X}} \{dX + \min_{Y \geq 0} \{eY \mid BY \geq f - AX\}\} \quad (19)$$

$$= \min_{X \in \mathcal{X}} \{dX + \max_{\lambda \geq 0} \{(f - AX)\lambda \mid B'\lambda \leq e\}\}, \quad (20)$$

where λ denotes the dual variables for the subproblem with only Y . The inner subproblem is called the *slave*, while the outer problem with only X is the *master*. Eq. (19) is via decomposition. Eq. (20) is via strong duality of the slave as an LP. Note in (20), the dual slave's feasibility region is independent of X . Let $\Lambda = \{\lambda \geq 0 \mid B'\lambda \leq e\}$ be the dual feasibility region of the slave. We can re-write (20) using only the extreme points Λ_o^* and extreme rays Λ_f^* in Λ :

$$\min_{X \in \mathcal{X}, \theta} \quad dX + \theta \quad (21a)$$

$$\text{s.t.} \quad \theta \geq (f - AX)\lambda_o, \quad \forall \lambda_o \in \Lambda_o^*, \quad (21b)$$

$$(f - AX)\lambda_f \leq 0, \quad \forall \lambda_f \in \Lambda_f^*. \quad (21c)$$

An extreme ray λ_f is a direction to which the dual slave objective can increase without bound, which exists only if the primal slave is infeasible for some $X \in \mathcal{X}$. Eq. (21c) eliminates such an infeasible X by bounding such a ray to be below 0.

It is impractical to enumerate all extreme points and extreme rays for a given problem. Hence, Benders decomposition works on a relaxed version of (21) that includes only a subset of the extreme points $\Lambda_o \subseteq \Lambda_o^*$ and extreme rays $\Lambda_f \subseteq \Lambda_f^*$:

$$\min_{X \in \mathcal{X}, \theta} \quad dX + \theta \quad (22a)$$

$$\text{s.t.} \quad \theta \geq (f - AX)\lambda_o, \quad \forall \lambda_o \in \Lambda_o, \quad (22b)$$

$$(f - AX)\lambda_f \leq 0, \quad \forall \lambda_f \in \Lambda_f. \quad (22c)$$

If a solution (X, θ) obtained by solving (22) fully satisfies (21b) and (21c), then (X, θ) is optimal to the original problem. This can be checked by solving the dual slave subproblem $\bar{\theta} = \max_{\lambda \geq 0} \{(f - AX)\lambda \mid B'\lambda \leq e\}$ and checking whether $\theta \geq \bar{\theta}$. If $\bar{\theta} = \infty$, then X is infeasible, and a new *feasibility cut* $(f - AX)\bar{\lambda} \leq 0$ can be added, where $\bar{\lambda}$ is an extreme ray of the unbounded dual. Otherwise if $\theta < \bar{\theta}$, an *optimality cut* $\theta \geq (f - AX)\bar{\lambda}$ can be added, where $\bar{\lambda}$ is the optimal dual solution. Following this intuition, the Benders decomposition starts with the relaxed problem (22) with no cuts, and then gradually solves the relaxed problem and the corresponding dual slave to progressively add cuts to (22). In any iteration, the value $dX + \theta$ obtained by solving (22) is a lower bound of the optimal value (since the problem is relaxed), while any $dX + \bar{\theta}$ where $\bar{\theta} \leq \infty$ obtained by solving the dual slave is an upper bound (since it

represents a feasible solution of the problem). The algorithm stops either when $\theta \geq \bar{\theta}$ which yields an optimal solution, or when (22) is infeasible which means the original problem is infeasible as well.

4.3 Nested Benders Decomposition

The original Benders decomposition works for two-stage problems. For our problem, we need the *nested Benders decomposition* [6]. As the name suggests, the algorithm stacks two layers of the original Benders decomposition: the *Phase-I algorithm* for SD–NPR, and the *Phase-II algorithm* for NPR–NE. As the middle layer, NPR is both the *slave* in Phase-I, and the *master* in Phase-II. In Phase-I, each iteration solves the relaxed SD denoted as O^{re} , and then solves the resulting $\mathcal{P}_t(X)$ for each t through the Phase-II; a feasibility cut as in (22c) or a set of optimality cuts as in (22b) is added to O^{re} at the end of a Phase-I iteration. In Phase-II, each iteration solves the relaxed NPR denoted as $\mathcal{P}_t^{\text{re}}(X)$, and then solves the resulting $Q_{t,k}(Y, R)$ for each sample k ; a feasibility cut or an optimality cut is added to $\mathcal{P}_t^{\text{re}}(X)$ at the end of a Phase-II iteration. After Phase-II converges, the final $\mathcal{P}_t^{\text{re}}(X)$ is equivalent to $\mathcal{P}_t(X)$, and hence the optimal dual solution of $\mathcal{P}_t^{\text{re}}(X)$ is used to construct the Phase-I optimality cuts. If Phase-II is infeasible for some X , then the dual (which has an extreme ray) is used to construct the Phase-I feasibility cut. The process stops when Phase-I converges or is deemed infeasible.

Below we formally derive the Benders cuts, starting from Phase-II. Given fixed \bar{X} , let $(\bar{Y}_t, \bar{r}(t))$ be an arbitrary Phase-II master solution. Assume $Q_{t,k}(\bar{Y}, \bar{R})$ is feasible for $\forall k$. Let $\bar{\pi} = \{\bar{\pi}_{t,k}^{(5)}(p)\} \cup \{\bar{\pi}_{t,k}^{(6)}(a)\} \cup \{\bar{\pi}_{t,k}^{(14c)}\}$ be the optimal dual multipliers for Constraints (5), (6) and (14c) respectively. Define q as the auxiliary variable in $\mathcal{P}_t^{\text{re}}(\bar{X})$ as θ in (22). The optimality cut to be added to $\mathcal{P}_t^{\text{re}}(\bar{X})$ is

$$q \geq \sum_k \mathcal{L}_{t,k}^{\mathcal{P}}(\bar{\pi}, Y, R) \quad (23)$$

where

$$\mathcal{L}_{t,k}^{\mathcal{P}}(\pi, Y, R) = \sum_{a \in A} \delta_{t,a}^k \pi_{t,k}^{(6)}(a) - \sum_{p \in P} \pi_{t,k}^{(5)}(p) y(t, p) - \pi_{t,k}^{(14c)} r(t).$$

If for some k , $Q_{t,k}(\bar{Y}, \bar{R})$ is infeasible. Let $\tilde{\pi}$ be an extreme ray of the dual of $Q_{t,k}(\bar{Y}, \bar{R})$. The feasibility cut to be added to $\mathcal{P}_t^{\text{re}}(\bar{X})$ is

$$\mathcal{L}_{t,k}^{\mathcal{P}}(\tilde{\pi}, Y, R) \leq 0. \quad (24)$$

Let dual solutions $\bar{\Pi}$ and extreme rays $\tilde{\Pi}$ be involved in $\mathcal{P}_t^{\text{re}}(\bar{X})$ in a Phase-II iteration. Then $\mathcal{P}_t^{\text{re}}(\bar{X}) \triangleq$

$$\begin{aligned} & \min_{Y_t \geq 0, r(t), q} \quad \rho_2 r(t) + q/K \\ & \text{s.t.} \quad (3)-(4); (23) \quad \forall \bar{\pi}_i \in \bar{\Pi}; (24) \quad \forall \tilde{\pi}_j \in \tilde{\Pi}. \end{aligned} \quad (25)$$

Now, given \bar{X} , assume $\mathcal{P}_t^{\text{re}}(\bar{X})$ is feasible for $\forall t$. Let $\bar{\gamma} = \{\bar{\gamma}_t^{(3)}(p)\} \cup \{\bar{\gamma}_t^{(4)}(l)\} \cup \{\bar{\gamma}_t^{(23)}(i)\} \cup \{\bar{\gamma}_t^{(24)}(j)\}$ be the optimal dual multipliers for (3), (4), (23) and (24) respectively. Define

$$\begin{aligned} \mathcal{L}_t^O(\gamma, X) = & \sum_{\bar{\pi}_i \in \bar{\Pi}} \bar{\gamma}_t^{(23)}(i) \sum_k \sum_{a \in A} \delta_{t,a}^k \bar{\pi}_{t,k}^{(6)}(a) + \\ & \sum_{\tilde{\pi}_j \in \tilde{\Pi}} \bar{\gamma}_t^{(24)}(j) \sum_{a \in A} \delta_{t,a}^{k_j} \tilde{\pi}_{t,k_j}^{(6)}(a) - \\ & \left(\sum_{p \in P} b_p^{\max} \bar{\gamma}_t^{(3)}(p) x(h_p) + \sum_{l \in L} b_l \bar{\gamma}_t^{(4)}(l) \right), \end{aligned}$$

where k_j is the infeasible sample corresponding to the j -th extreme ray added to $\mathcal{P}_t^{\text{re}}(\bar{X})$. The optimality cuts to O^{re} (Phase-I) is

$$D \geq \mathcal{L}_t^O(\bar{\gamma}, X), \quad \forall t = 1 \dots T. \quad (26)$$

Note in Phase-I, variable D serves exactly as the auxiliary variable θ in O^{re} . Instead of adding one optimality cut per iteration, T cuts are added in the form of (26). Later in Sec. 4.4, we show that this idea can be extended to Phase-II to expedite algorithm convergence. Now suppose for some t , $\mathcal{P}_t^{\text{re}}(\bar{X})$ is infeasible. Let \bar{y} be an extreme ray of the dual of $\mathcal{P}_t^{\text{re}}(\bar{X})$. The feasibility cut to O^{re} is

$$\mathcal{L}_t^O(\bar{y}, X) \leq 0. \quad (27)$$

We finally address convergence. For Phase-II, let $Q_{t,k}^*$ be the optimal value of $Q_{t,k}(\bar{Y}, \bar{R})$, then an *upper bound* of $\mathcal{P}_t(\bar{X})$ is $\rho_2 \bar{r}(t) + \sum_k Q_{t,k}^*$ (since $\{Q_{t,k}^* \mid k\}$ represents a feasible solution to NE); let \mathcal{P}_t^* be the optimal value of $\mathcal{P}_t^{\text{re}}(\bar{X})$, then \mathcal{P}_t^* is a *lower bound* of $\mathcal{P}_t(\bar{X})$ (since it is the optimal value of the *relaxed* problem). Phase-II convergence thus happens when $\mathcal{P}_t^* \geq \rho_2 \bar{r}(t) + \sum_k Q_{t,k}^*$. Similarly for Phase-I, an *upper bound* of O is $\max_t \mathcal{P}_t^*$ at Phase-II convergence; a *lower bound* is the optimal value O^* of O^{re} . Phase-I convergence happens when $O^* \geq \max_t \mathcal{P}_t^*$. Note that for either phase, the lower bound is non-decreasing through iterations (since the master becomes “less” relaxed in each iteration), yet the upper bound may not be non-increasing (since different feasible solutions may be found with various values). The algorithm thus needs to track the best upper bound ever found in all iterations to verify convergence.

Algorithm 1 shows the full algorithm. Lines 2–30 show Phase-I algorithm for solving O , while Lines 7–24 show Phase-II algorithm for solving $\mathcal{P}_t(\bar{X})$ for each t . Note that for middle-layer $\mathcal{P}_t^{\text{re}}(\bar{X})$, the Phase-II cuts sc are shared across Phase-I iterations. This is because although a Phase-II cut generated in a previous Phase-I iteration may not be a valid cutting plane for the current $\mathcal{P}_t(\bar{X})$, it is a valid cut for the overall SD–NPR problem. Sharing cuts across Phase-I iterations can prune suboptimal X solutions without running the full Phase-II on them, and hence can speed-up convergence.

Algorithm 1 exploits the horizontal decomposition in both Phase-I and Phase-II, where $\mathcal{P}_t^{\text{re}}(\bar{X})$ and $Q_{t,k}(\bar{Y}, \bar{R})$ subproblems are solved independently for each t and/or k respectively. This is enabled by separating all *complicating variables and constraints* into the corresponding master problems, *i.e.*, variables (X, D) and constraints (1)–(2) for SD–NPR, and $(Y, r(t))$ and (3)–(4) for NPR–NE. Remember we separated variables R and W when defining $\mathcal{P}_t(X)$ and $Q_{t,k}(Y, R)$ to enable such a decomposition. With the fastest LP algorithm running in no less than cubic time [35], such decompositions can reduce the dependence of time complexity over T and K from cubic to linear for solving the subproblems. which is specifically important when the sample size K is large. This is our main motivation for applying Benders decomposition.

Since (14) is an MILP with finite master feasibility set and LPs as subproblems, we state Algorithm 1’s convergence as below, whose proof is established in [3, 6] and hence omitted due to page limit.

THEOREM 4.1. *Algorithm 1 converges to the optimal solution of a feasible REP instance or returns “Infeasible” for an infeasible instance in finite iterations.* \square

4.4 Efficiency Enhancement Techniques

Though Algorithm 1 converges, it may have a slow convergence as a cutting plane method. Below, we describe three computational techniques to improve its efficiency. We start with a novel analytical dual solving technique for the NE subproblems. We then describe

Algorithm 1: Nested Benders Decomposition for REP

Input: Network G , pathbook P , cost budget C , samples $\{S_t\}$, convergence tolerance ϵ .

Output: Optimal (X, Y, Z, D, R, W) .

```

1   $LB^m \leftarrow -\infty, UB^m \leftarrow \infty, mc \leftarrow \emptyset, sc \leftarrow \emptyset;$ 
2  repeat // Phase-I iteration
3      Solve  $LB^m = O^{\text{re}}$  with cuts  $mc$  for  $(\bar{X}, \bar{D})$ ;
4      if  $O^{\text{re}}$  infeasible then return Infeasible;
5      for  $t = 1 \dots T$  do
6           $LB_t^s \leftarrow -\infty, UB_t^s \leftarrow \infty;$ 
7          repeat // Phase-II iteration
8              Solve  $LB_t^s = \mathcal{P}_t^{\text{re}}(\bar{X})$  with cuts  $sc$ , and denote the
              solution as  $(\bar{Y}_t, \bar{r}(t))$ ;
9              if  $\mathcal{P}_t^{\text{re}}(\bar{X})$  infeasible then
10                  Get extreme ray of dual  $\mathcal{P}_t^{\text{re}}(\bar{X})$ ;
11                  Add Phase-I feasibility cut to  $mc$ ;
12                  break;
13              for  $k = 1 \dots K$  do
14                  Solve  $Q_{t,k}^* = Q_{t,k}(\bar{Y}, \bar{R})$  for  $(\bar{Z}, \bar{W})$ ;
15                  if  $Q_{t,k}(\bar{Y}, \bar{R})$  feasible then
16                      Record optimal dual of  $Q_{t,k}(\bar{Y}, \bar{R})$ ;
17                  else
18                      Get extreme ray of dual  $Q_{t,k}(\bar{Y}, \bar{R})$ ;
19                      Add Phase-II feasibility cut to  $sc$ ;
20                      break;
21              if  $\forall k, Q_{t,k}(\bar{Y}, \bar{R})$  is feasible then
22                   $UB_t^s = \min\{UB_t^s, \rho_2 \bar{r}(t) + \sum_k Q_{t,k}^*\};$ 
23                  Construct Phase-II optimality cut from
                   $Q_{t,k}(\bar{Y}, \bar{R})$  duals, and add cut to  $sc$ ;
24          until  $UB_t^s - LB_t^s \leq \epsilon$ ;
25          if  $\mathcal{P}_t^{\text{re}}(\bar{X})$  infeasible then break;
26          else Record optimal dual of  $\mathcal{P}_t^{\text{re}}(\bar{X})$ ;
27      if  $\forall t, \mathcal{P}_t^{\text{re}}(\bar{X})$  is feasible then
28           $UB^m = \min\{UB^m, \max_t\{LB_t^s\}\};$ 
29          Construct Phase-I optimality cut from recorded
           $\mathcal{P}_t^{\text{re}}(\bar{X})$  duals, and add cut to  $mc$ ;
30  until  $UB^m - LB^m \leq \epsilon$ ;
31  return  $\bar{X}$  and best recorded subproblem solutions.
```

two established techniques, one applying to the original Benders decomposition, and one applying to the nested algorithm.

Analytical Solving for Phase-II Slave Dual: This technique applies due to the structure of the NE subproblem, *i.e.*, finding optimal per-sample demand allocation across paths is easy. Given a sample k , for each $a \in A$, one can enumerate its paths towards all hosts in increasing delays, and then allocate demand to the shortest paths until all demand is allocated. This leads to a linear-time optimal algorithm for primal NE, assuming pre-sorted paths in all samples.

Yet, dual information is needed to add cut(s). In general, obtaining the dual of a primal solution is equivalent to solving a system of linear equations, which takes cubic time, not much faster than solving the dual LP directly. For our problem, though, we propose a *linear time analytical solution* (with pre-sorting) for dual NE.

Consider an instance of $Q_{t,k}(\bar{Y}, \bar{R})$. Let $\pi = \{\pi_{t,k}^{(5)}(p)\} \cup \{\pi_{t,k}^{(6)}(a)\} \cup \{\pi_{t,k}^{(14c)}\}$ be the dual variables for Constraints (5), (6) and (14c) respectively. The dual of $Q_{t,k}(\bar{Y}, \bar{R})$ is

$$\max_{\pi \geq 0} \sum_{a \in A} \left(\delta_{t,a}^k \pi_{t,k}^{(6)}(a) - \sum_{p \in P_a} \bar{y}(t,p) \pi_{t,k}^{(5)}(p) \right) - \bar{r}(t) \pi_{t,k}^{(14c)} \quad (28a)$$

$$\text{s.t. } \pi_{t,k}^{(14c)} \leq \frac{\rho_2}{1-\alpha}; \quad (28b)$$

$$\pi_{t,k}^{(6)}(a) - \pi_{t,k}^{(5)}(p) \leq \frac{1}{\delta_t^k} d_{t,p}^k \left(\rho_1 + \pi_{t,k}^{(14c)} \right), \quad (28c)$$

$$\forall a \in A, p \in P_a.$$

Now consider an AP $a \in A$. Let its path set P_a be sorted as $\mathbf{P}_a = (p_1^a, p_2^a, \dots, p_p^a)$ in increasing order of path delays $d_{t,p}^k$. For simplicity we assume all path delays are distinct. Define the *critical path index* as $j_a = \arg \min_j \{d_{t,p_j^a}^k \mid p_j^a \in P_a, \sum_{i=1}^j \bar{y}(t, p_i^a) \geq \delta_{t,a}^k\}$ and $d_{\text{crit}}^a = d_{t,p_{j_a}^a}^k$, i.e., the min-delay path $p_{j_a}^a$ such that all paths with delay no greater than $d_{t,p_{j_a}^a}^k$ can fully satisfy the demand $\delta_{t,a}^k$. Based on the aforementioned optimal primal solution, any path with $d_{t,p}^k < d_{\text{crit}}^a$ will have a saturated allocation $z(t,k,p) = \bar{y}(t,p)$, while all paths with $d_{t,p}^k > d_{\text{crit}}^a$ will have $z(t,k,p) = 0$.

Further define $\bar{y}_{\text{sum}}^a = \sum_{j < j_a} \bar{y}(t, p_j^a)$ and $\bar{y}_{\text{crit}}^a = \delta_{t,a}^k - \bar{y}_{\text{sum}}^a$. From the primal solution, the optimal average delay can be computed as $D_{t,k}^* = \frac{1}{\delta_t^k} \sum_a \left(\sum_{j < j_a} d_{t,p_j^a}^k \bar{y}(t, p_j^a) + d_{\text{crit}}^a \bar{y}_{\text{crit}}^a \right)$. We then compute an optimal dual solution by the following:

$$\pi_{t,k}^{(14c)} = \begin{cases} \frac{\rho_2}{1-\alpha} & D_{t,k}^* \geq \bar{r}(t) \\ 0 & \text{otherwise} \end{cases}, \quad (29a)$$

$$\pi_{t,k}^{(6)}(a) = \frac{1}{\delta_t^k} \left(\rho_1 + \pi_{t,k}^{(14c)} \right) d_{\text{crit}}^a, \quad (29b)$$

$$\pi_{t,k}^{(5)}(p) = \frac{1}{\delta_t^k} \left(\rho_1 + \pi_{t,k}^{(14c)} \right) \left(d_{\text{crit}}^a - d_{t,p}^k \right)^+. \quad (29c)$$

To compute these values, the paths need to be sorted, which takes $O(p \log p)$ time. The sorting only needs to be performed once for every sample. Therefore, with a $O(TKP \log P)$ pre-sorting, each dual of $Q_{t,k}(\bar{Y}, \bar{R})$ can be analytically solved in linear time using Eqs. (29). We note that feasibility of $Q_{t,k}(\bar{Y}, \bar{R})$ can be addressed during the computation: if $\sum_{p \in P_a} \bar{y}(t,p) \geq \delta_{t,a}^k$ for $\forall a \in A$, then $Q_{t,k}(\bar{Y}, \bar{R})$ is feasible. Otherwise, (29) also presents a valid extreme ray when d_{crit}^a takes any value that is greater than the longest path delay with positive $\bar{y}(t,p)$ and $\pi_{t,k}^{(14c)} = \frac{\rho_2}{1-\alpha}$; the proof is similar to that of Theorem 4.2 and hence is omitted. An alternative for addressing infeasible $Q_{t,k}(\bar{Y}, \bar{R})$ is to add a new constraint

$$\sum_{p \in P_a} y(t,p) \geq \max_k \{\delta_{t,p}^k\}, \quad \forall a \in A$$

to the Phase-II master $\mathcal{P}_t^{\text{re}}(X)$, and modify the cuts in (26) and (27). This ensures that $Q_{t,k}(\bar{Y}, \bar{R})$ is always feasible. The following theorem states the optimality of Eqs. (29) when the dual NE subproblem is feasible given \bar{Y} and \bar{R} :

THEOREM 4.2. *For a feasible instance of $Q_{t,k}(\bar{Y}, \bar{R})$, Eqs. (29) denote an optimal solution to Program (28). \square*

PROOF. Let us consider a specific $a \in A$ and its path set P_a . To start with, we fix $\pi_{t,k}^{(14c)}$ and define $\eta_{t,k} = \frac{1}{\delta_t^k} \left(\rho_1 + \pi_{t,k}^{(14c)} \right)$. For a given $\pi_{t,k}^{(6)}(a)$, divide P_a into four sets:

- $P_a^1 = \{p \mid \bar{y}(t,p) > 0, \eta_{t,k} \cdot d_{t,p}^k \leq \pi_{t,k}^{(6)}(a)\}$,
- $P_a^2 = \{p \mid \bar{y}(t,p) > 0, \eta_{t,k} \cdot d_{t,p}^k > \pi_{t,k}^{(6)}(a)\}$,
- $P_a^3 = \{p \mid \bar{y}(t,p) = 0, \eta_{t,k} \cdot d_{t,p}^k \leq \pi_{t,k}^{(6)}(a)\}$, and
- $P_a^4 = \{p \mid \bar{y}(t,p) = 0, \eta_{t,k} \cdot d_{t,p}^k > \pi_{t,k}^{(6)}(a)\}$.

For $p \in P_a^2, \pi_{t,k}^{(5)}(p) = 0$ by complementary slackness of Constraint (5). For $p \in P_a^3 \cup P_a^4$, $\pi_{t,k}^{(5)}(p)$ takes arbitrary value no less than $\max\{\pi_{t,k}^{(6)}(a) - \eta_{t,k} d_{t,p}^k, 0\}$, which ensures (28c); note that none of these paths contribute to the objective value since $\bar{y}(t,p) = 0$. For $p \in P_a^1$, we must have a strict equality:

$$\pi_{t,k}^{(5)}(p) = \pi_{t,k}^{(6)}(a) - \eta_{t,k} d_{t,p}^k, \quad \forall p \in P_a^1. \quad (30)$$

If any such $\pi_{t,k}^{(5)}(p)$ is more, then we can solely decrease $\pi_{t,k}^{(5)}(p)$ to increase the objective value, contradicting optimality.

Taking (30) into the objective function (28a), we can re-write the first part of (28a) for each $a \in A$ as:

$$\mathcal{D}_t^k(a) = \delta_{t,a}^k \pi_{t,k}^{(6)}(a) - \sum_{p \in P_a} \bar{y}(t,p) \pi_{t,k}^{(5)}(p) \quad (31a)$$

$$= \delta_{t,a}^k \pi_{t,k}^{(6)}(a) - \sum_{p \in P_a^1} \bar{y}(t,p) \left(\pi_{t,k}^{(6)}(a) - \eta_{t,k} d_{t,p}^k \right) \quad (31b)$$

$$= \left(\delta_{t,a}^k - \sum_{p \in P_a^1} \bar{y}(t,p) \right) \pi_{t,k}^{(6)}(a) + \sum_{p \in P_a^1} \bar{y}(t,p) \eta_{t,k} d_{t,p}^k. \quad (31c)$$

In (31c), the right term is only related to the set P_a^1 , while the left term is decided by both P_a^1 and $\pi_{t,k}^{(6)}(a)$. Recall we define the *critical delay* as $d_{\text{crit}}^a = \min\{d \mid \sum_{p: d_{t,p}^k \leq d} \bar{y}(t,p) \geq \delta_{t,a}^k\}$ (equivalent to the definition by a critical path). If $\pi_{t,k}^{(6)}(a) > \eta_{t,k} d_{\text{crit}}^a$, by definition $\delta_{t,a}^k - \sum_{p \in P_a^1} \bar{y}(t,p) \leq 0$. Then we can reduce $\pi_{t,k}^{(6)}(a)$ until $\pi_{t,k}^{(6)}(a) = \eta_{t,k} d_{\text{crit}}^a$ to increase (or not to decrease if $\delta_{t,a}^k = \sum_{p \in P_a^1} \bar{y}(t,p)$) the value of (31c). If $\pi_{t,k}^{(6)}(a) < \eta_{t,k} d_{\text{crit}}^a$, by definition $\delta_{t,a}^k - \sum_{p \in P_a^1} \bar{y}(t,p) > 0$. Then we can increase $\pi_{t,k}^{(6)}(a)$ until $\pi_{t,k}^{(6)}(a) = \eta_{t,k} d_{\text{crit}}^a$ to again increase (31c). The only flexibility $\pi_{t,k}^{(6)}(a)$ can have is when $\delta_{t,a}^k = \sum_{p: d_{t,p}^k \leq d_{\text{crit}}^a} \bar{y}(t,p)$, in which case $\pi_{t,k}^{(6)}(a)$ can take any value between $\eta_{t,k} d_{\text{crit}}^a$ and $\eta_{t,k} d_{t,p_n}^k$ where p_n is the next path after the critical path that has a positive $\bar{y}(t,p)$.

The above has addressed the optimal values of $\pi_{t,k}^{(5)}(p)$ and $\pi_{t,k}^{(6)}(a)$ as in (29b) and (29c). Now if we take these values and $\eta_{t,k}$ into (31c), then $\sum_a \mathcal{D}_t^k(a) = D_{t,k}^* \cdot (\rho_1 + \pi_{t,k}^{(14c)})$ where $D_{t,k}^*$ is the optimal sample average delay. The full dual objective can be re-written as

$$\mathcal{D}_t^k = \sum_{a \in A} \mathcal{D}_t^k(a) - \bar{r}(t) \pi_{t,k}^{(14c)} \quad (32a)$$

$$= \rho_1 \cdot D_{t,k}^* + (D_{t,k}^* - \bar{r}(t)) \pi_{t,k}^{(14c)}. \quad (32b)$$

It is now clear that if $D_{t,k}^* - \bar{r}(t) \geq 0$, then $\pi_{t,k}^{(14c)}$ can take the maximum value $\frac{\rho_2}{1-\alpha}$ constrained by (28b); otherwise $\pi_{t,k}^{(14c)} = 0$. \square

Multiple Cuts: In each Phase-II iteration in Algorithm 1, at most one optimality cut (in the form of (22b)) is added. Birge and Louveaux [7] first noticed that by dividing a single optimality cut into multiple cuts, one for each independent subproblem for example, the convergence speed can be improved. Intuitively, this is because multiple cuts can prune more sub-optimal region of the master solution space in each iteration, hence the number of iterations is reduced. We can apply this technique to our algorithm. Instead of having one auxiliary variable q in Eq. (23), K variables are defined as $\{q(k)\}$. The optimality cuts to replace Eq. (23) are

$$q(k) \geq \mathcal{L}_{t,k}^{\mathcal{P}}(\bar{\pi}, Y, R), \quad \forall k, \quad (33)$$

and the objective function of $\mathcal{P}_t^{\text{re}}(\bar{X})$ in (25) is updated as

$$\min \rho_2 r(t) + \frac{1}{K} \sum_k q(k). \quad (34)$$

The Phase-I algorithm has already employed a multi-cut formulation as in (26). We highlight that this technique does not change the execution in Algorithm 1: all subproblems $Q_{t,k}(\bar{Y}, \bar{R})$ must still be solved before the cuts in (33) can be added, in order to avoid adding optimality cuts for an infeasible (\bar{Y}, \bar{R}) . This is specifically important when implementing the algorithm distributedly. We show in Sec. 5 that this can achieve significant speed-up of the algorithm.

Fast Forward Fast Backward (FFFB): In Algorithm 1, each Phase-I iteration consists of a full Phase-II convergence: only dual of the optimal Phase-II master solution is passed onto Phase-I. This is the *Fast Forward (FF)* approach proposed by Birge [5]. Another choice is *Fast Backward (FB)*: its main loop consists of the Phase-II algorithm, where each Phase-II iteration consists of a full convergence of Phase-I [31]. However, neither approach achieves the best practical performance. A balanced approach, called *Fast Forward Fast Backward (FFFB)*, is shown to be the most efficient in practice [6]. In each main iteration, the algorithm performs a full forward sweep (solving Phase-I master, Phase-II master and Phase-II slave), followed by a full backward sweep (adding Phase-II cut(s), solving Phase-II master *once*, and then adding Phase-I cut(s) based on it); if in the forward sweep the Phase-II master is infeasible, then such information is directly passed to the Phase-I master. FFB speeds-up convergence similarly as sharing Phase-II cuts: cuts based on non-optimal (but feasible) NPR help eliminate suboptimal Phase-I master solutions early, which reduces the total number of Phase-I master solutions visited. We show its effectiveness also in Sec. 5.

4.5 Discussions

Complexity. Algorithm 1 converges to optimal (Theorem 4.1) but has exponential worst-case complexity, and techniques in Sec. 4.4 do not reduce it to polynomial. However, the techniques (especially dual analytical solving) present two main advantages in practice. First, we found that practically NE solving is the efficiency bottleneck due to the large number of samples needed for accurate sampling. Our technique speeds-up worst-case per-iteration NE solving by at least two orders of magnitudes, greatly improving efficiency without sacrificing solution quality. Second, iterative algorithms are commonly implemented with a time bound in practice. Faster per-iteration solving greatly increases the solution space searched in bounded time, thus improving solution quality.

Distributed computation. Algorithm 1 can be distributed among ECM, NM and ESP. Each party can solve its own subproblem (SD

for ECM, NPR for NM, and NE for ESP itself), and exchange solution information. Distributed computation may not only reduce per-party computational overhead, but also preserves privacy of certain parties, e.g., the computing resource distribution of the ECM, and the network topology of the NM. We will explore distributed computation and its implications in our future work.

5 PERFORMANCE EVALUATION

5.1 Evaluation Method

We synthesized real-world demand traces with randomly generated data for evaluation. For demands, we used the NYC 2018 Yellow Taxi Trip Data [27], which contains around 112 million taxi trips including drop-off times and taxi zones. We divided the year into days as periods, and each day was divided into 2–24 (4 by default) time slots. We regarded each taxi zone as a single AP, and took the average number of passengers dropped-off in each slot as a demand sample at the AP, where each passenger represented 1 Mbps of fixed demand. We picked the top 5 (*small*) or 20 (*medium*) most popular drop-off zones over the entire 262 zones depending on the desired problem size, which accounted for 18% or 55% of all demands respectively. We used the first 100 days as our *training samples*, i.e., samples used to solve the optimization problem to derive SD and NPR. We then used the next 265 days as our *test samples*, i.e., samples used to evaluate our solution once it is employed.

We synthesized the other inputs of the problem. First, we built random topologies among APs using the Watts-Strogatz model [30] with $k = 4$ and $p = 0.3$. For the small-sized problem with 5 APs, we added 20 artificial network nodes to construct the topology. We then randomly picked 5 nodes, each connected to an edge computing node through an infinite-capacity link. Deployment costs were generated in a normal distribution $\mathcal{N}(\mu, \sigma^2)$ with $\mu = 1000$ and $\sigma = 200$, while the ESP's cost budget was 3300. For the pathbook, we ran the k -shortest path algorithm to generate 3 min-hop paths for each AP-host pair. Such a number has been reported in [16] to be able to admit over 90% of the maximum flow traffic in a typical network. We simulated two scenarios: a *normal* scenario where link capacities (except for links to edge nodes) were fixed at 5 Gbps and delays were generated in $\mathcal{N}(10, 4)$ ms, and a *congested* scenario where capacities were downgraded to 2 Gbps and half edge nodes randomly selected in each time slot had 50× prolonged connection delays (simulating about half packets experiencing 1s retransmission timeouts due to host congestion).

By default, we set $\rho_1 = \rho_2 = 0.5$ and $\alpha = 0.95$. For convergence, we set $\epsilon = 10^{-3}$. To reduce random noise, we averaged results over 10 independent runs on inputs generated with different seeds. We implemented all algorithms in Python, where LPs/MILPs were solved using MOSEK [2]. Running time-related simulations were run on a Windows PC with Hex-Core 3.2GHz CPU and 16GB memory.

5.2 Experiment Results

5.2.1 Mean vs. CVaR. Fig. 2 shows the trade-off between expectation and CVaR with α and (ρ_1, ρ_2) for the training set. We conducted simulations on the *medium* problem size and with *normal* load. In Fig. 2(a), we varied α from 0.01 to 0.99. With higher confidence α , the CVaR that measures the worst-case delay beyond such a confidence is higher. When α is close to 0, the CVaR converges

to the expectation; when α is close to 1, the CVaR approaches the worst-case performance observed among all samples. However, optimizing the solution with a higher α has little impact on the expected value. In Fig. 2(b), by changing ρ_1 (with fixed ρ_2), we can optimize different combinations of expectation and CVaR. A higher ρ_1 results in a lower expected value but a higher CVaR. In practice, ESP can fine-tune these parameters to find its best trade-off.

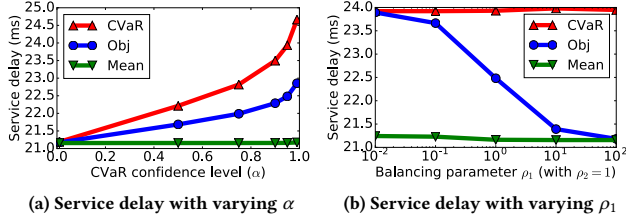


Figure 2: Trade-off between sample mean and CVaR

5.2.2 Time-Varying vs. Time-Agnostic. Fig. 3 compares two modeling choices: time-varying (**T-Var**) and time-agnostic (**T-Ago**). Both models were solved using Algorithm 1. In **T-Ago**, all samples were flattened out across time slots, and the problem was solved as if there was just one time slot. We ran our simulations in the *small* setting with *congested* load, and varied the number of time slots from 2 to 24. Remember we used average demands in a time slot as a sample. With more slots, the input includes more samples (representing more dynamics), and hence both models lead to higher CVaR and thus higher objective value. **T-Var** has a growing performance advantage over **T-Ago**, because of the former's ability to adjust to different demand and network patterns at different times and perform finer-grained optimization with more slots.

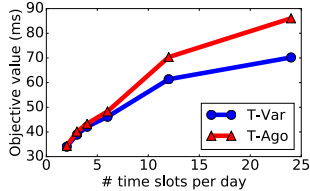


Figure 3: Objective with time-varying or fixed provisioning

5.2.3 Optimal vs. Heuristic. Fig. 4 shows the training and testing objective values obtained with different host numbers, where we compared our algorithm to two heuristics in terms of the quality of solution (objective value). Here, **REP** refers to our algorithm; **RAND** refers to a random host selection heuristic with NPR-NE solved using our Phase-II; **AVG** is a heuristic that simply optimizes the average delay for each time slot when delays and demands are simply averaged over all samples, which is solved in one integrated MILP. We ran our simulations in the *medium* setting with *normal* load. With more hosts, the objectives of **REP** and **AVG** are generally decreasing. **RAND**, however, has fluctuating objective values due to random deployment. Among the algorithms, **REP** naturally achieves the lowest objective value due to its optimality, though **AVG** has performance not significantly worse than ours. However, we emphasize that our advantageous solution will constantly impact service performance in the long run after deployment, based on the test sample results in Fig. 4(b). Hence it is worthwhile to improve performance using our algorithm in initial provisioning.

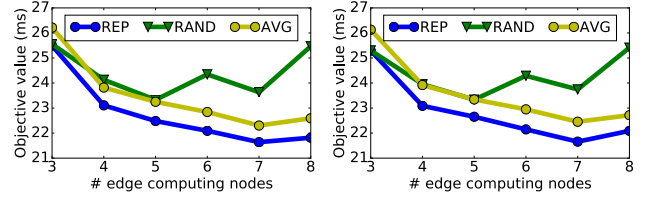


Figure 4: Comparison of our algorithm with heuristics

5.2.4 Efficiency Enhancement. Finally, we evaluate the efficiency enhancement techniques in Fig. 5. We compared **Full** (our algorithm with all three techniques), to **!Multi** (**Full** without multi-cut), **!FFFB** (**Full** without FFFB), and **!Ana** (**Full** without the analytical dual solution and solves Phase-II slave duals through direct LP solving). We ran our simulations in the *small* setting with *normal* load. As we can see, each technique achieves multifold improvement of convergence speed. Since these techniques are orthogonal to each other, the overall speed-up is in two orders of magnitude, which demonstrates the effectiveness of these techniques.

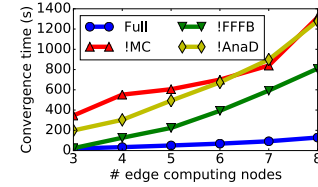


Figure 5: Convergence with/without efficiency techniques

6 RELATED WORK

6.1 Resource Allocation in Edge Computing

There has been two categories of efforts on edge resource allocation. One focuses on allocating non-network resources, stemming from existing research on distributed computing. Most such efforts consider fixed network allocations and statistics (e.g., delay). Examples include task scheduling [8, 26, 38], task offloading / service placement [10, 28, 29, 32], resource allocation [9], caching [34], etc. The other recognizes criticality of the network in service performance, and employs active network planning to best utilize network information and/or allocate network resource. For instance, Gao *et al.* [12] explored how access networks selection can improve service QoS. Josilo *et al.* [13] and Liu *et al.* [19, 20] studied multi-dimensional resource allocation including both computing and network resources. Yu *et al.* [37] first identified the importance and challenge of complex edge network topologies, and proposed approximation schemes for QoS-guaranteed joint service deployment and network provisioning for real-time applications. Besides, Poularakis *et al.* [23], Ouyang *et al.* [22] and Meng *et al.* [21] also explored joint service placement and routing with storage, migration cost, and task deadline considerations, respectively.

This paper uses a similar service model as [37], where a service processes continuous data streams. However, this paper considers more practical factors, including deployment costs, time-varying and uncertain demands and network condition, robustness, etc. We note that robustness in edge resource allocation has yet been well explored, and has only been studied in some very recent work [18].

6.2 Other Related Work

Besides edge resource allocation, other related problems include virtual network or infrastructure embedding (VNE/VIE) [11] and service function chaining (SFC) [4]. Few efforts on these problems have explored the time-varying demands and network condition. Exceptions include Xie *et al.* [33] for VIE and Li *et al.* [17] for SFC. Xie *et al.* [33] extracted time-varying characteristics of peak application demands via profiling, and then made bandwidth allocation accordingly. Li *et al.* [17] used a time-varying model to compute correlations between service functions for joint deployment. Both regarded the time-varying demands as deterministic or estimatable, and did not consider robustness when facing uncertain dynamics.

Conventionally, CVaR is a risk measure widely used in economics and finance, notably in portfolio optimization [24]. However, it has recently found use in computing and networking for robust optimization as well. For example, Rullo *et al.* [25] and Yu *et al.* [36] both used CVaR to represent and optimize security risks in IoT.

7 CONCLUSIONS

In this paper, we studied robust edge resource provisioning with time-varying demands and network conditions. For an edge service, the service provider employs a three-stage approach, including service deployment, network provisioning, and performance estimation. We proposed a time-slotted stochastic model to capture the temporal-spatial demand distributions observed in real-world datasets. We introduced a widely used risk model in economics, Conditional Value-at-Risk, which enables trade-off between expected performance and robustness. Provisioning was formulated as a three-stage sample-based stochastic optimization problem. Based on the problem structure, we proposed a nested Benders decomposition approach to optimally solve the problem, and described several efficiency enhancement techniques to drastically speed-up convergence. We presented results from real dataset-based simulations, which demonstrated the advantages of the time-varying and robustness models, the proposed algorithm compared to heuristics, and the efficiency enhancement techniques, respectively.

ACKNOWLEDGMENTS

This research was supported in part by NSF grants 1704092, 1704662, 1717197, and 1717315, and JSPS KAKENHI Grant No. JP20H00592. The information reported here does not reflect the position or the policy of the funding agency.

REFERENCES

- [1] Ravindra K. Ahuja, Thomas L. Magnanti, and James B. Orlin. 1993. *Network Flows: Theory, Algorithms and Applications*.
- [2] MOSEK ApS. 2019. MOSEK Optimizer API for Python Release 9.0.98. <https://docs.mosek.com/9.0/pythonfusion/index.html>
- [3] Jacques F Benders. 1962. Partitioning Procedures for Solving Mixed-Variables Programming Problems. *Numer. Math.* 4, 1 (dec 1962), 238–252.
- [4] Deval Bhamare, Raj Jain, Mohammed Samaka, and Aiman Erbad. 2016. A Survey on Service Function Chaining. *J. Netw. Comput. Appl.* 75 (nov 2016), 138–155.
- [5] John R. Birge. 1985. Decomposition and Partitioning Methods for Multistage Stochastic Linear Programs. *Oper. Res.* 33, 5 (oct 1985), 989–1007.
- [6] John R. Birge and François Louveaux. 2011. *Introduction to Stochastic Programming*.
- [7] John R. Birge and François V. Louveaux. 1988. A Multicut Algorithm for Two-Stage Stochastic Linear Programs. *Eur. J. Oper. Res.* 34, 3 (mar 1988), 384–392.
- [8] Luiz F Bittencourt, Javier Diaz-Montes, Rajkumar Buyya, Omer F Rana, and Manish Parashar. 2017. Mobility-Aware Application Scheduling in Fog Computing. *IEEE Cloud Comput.* 4, 2 (mar 2017), 26–35.
- [9] Gabriele Castellano, Flavio Esposito, and Fulvio Risso. 2019. A Distributed Orchestration Algorithm for Edge Computing Resources with Guarantees. In *Proc. IEEE INFOCOM*. 2548–2556.
- [10] Ruilong Deng, Rongxing Lu, Chengzhe Lai, Tom Hao Luan, and Hao Liang. 2016. Optimal Workload Allocation in Fog-Cloud Computing Towards Balanced Delay and Power Consumption. *IEEE Internet Things J.* 3, 6 (2016), 1171–1181.
- [11] Andreas Fischer, Juan Felipe Botero, Michael Till Beck, Hermann de Meer, and Xavier Hesselbach. 2013. Virtual Network Embedding: A Survey. *IEEE Commun. Surv. Tutorials* 15, 4 (jan 2013), 1888–1906.
- [12] Bin Gao, Zhi Zhou, Fangming Liu, and Fei Xu. 2019. Winning at the Starting Line: Joint Network Selection and Service Placement for Mobile Edge Computing. In *Proc. IEEE INFOCOM*. 1459–1467.
- [13] Sladana Josilo and Gyorgy Dan. 2019. Wireless and Computing Resource Allocation for Selfish Computation Offloading in Edge Computing. In *Proc. IEEE INFOCOM*. 2467–2475.
- [14] Anton J. Kleywegt, Alexander Shapiro, and Tito Homem-de Mello. 2002. The Sample Average Approximation Method for Stochastic Discrete Optimization. *SIAM J. Optim.* 12, 2 (jan 2002), 479–502.
- [15] David Kotz, Tristan Henderson, Ilya Abyzov, and Jihwang Yeo. [n. d.]. CRAWDAD Dataset Dartmouth/Campus (V. 2009-09-09). <https://crawdad.org/dartmouth/campus/20090909>
- [16] Mathieu Leconte, Apostolos Destounis, and Georgios Paschos. 2018. Traffic Engineering with Precomputed Pathbooks. In *Proc. IEEE INFOCOM*. 234–242.
- [17] Defang Li, Peilin Hong, Kaiping Xue, and Jianing Pei. 2018. Virtual Network Function Placement Considering Resource Optimization and SFC Requests in Cloud Datacenter. *IEEE Trans. Parallel Distrib. Syst.* 29, 7 (jul 2018), 1664–1677.
- [18] Haijun Liao, Zhenyu Zhou, Shahid Mumtaz, and Jonathan Rodriguez. 2019. Robust Task Offloading for IoT Fog Computing Under Information Asymmetry and Information Uncertainty. In *Proc. IEEE ICC*. 1–6.
- [19] Qiang Liu and Tao Han. 2019. DIRECT: Distributed Cross-Domain Resource Orchestration in Cellular Edge Computing. In *Proc. ACM MobiHoc*.
- [20] Qiang Liu and Tao Han. 2019. VirtualEdge: Multi-Domain Resource Orchestration and Virtualization in Cellular Edge Computing. In *Proc. IEEE ICDCS*.
- [21] Jiaying Meng, Haisheng Tan, Chao Xu, Wanli Cao, Liuyan Liu, and Bojie Li. 2019. Dedas: Online Task Dispatching and Scheduling with Bandwidth Constraint in Edge Computing. In *Proc. IEEE INFOCOM*. 2287–2295.
- [22] Tao Ouyang, Rui Li, Xu Chen, Zhi Zhou, and Xin Tang. 2019. Adaptive User-managed Service Placement for Mobile Edge Computing: An Online Learning Approach. In *Proc. IEEE INFOCOM*. 1468–1476.
- [23] Konstantinos Poularakis, Jaime Llorca, Antonia M. Tulino, Ian Taylor, and Leandros Tassioulas. 2019. Joint Service Placement and Request Routing in Multi-cell Mobile Edge Computing Networks. In *Proc. IEEE INFOCOM*. 10–18.
- [24] R. Tyrrell Rockafellar and Stanislav Uryasev. 2000. Optimization of Conditional Value-at-Risk. *J. Risk* 2 (2000), 21–41.
- [25] Antonino Rullo, Edoardo Serra, Elisa Bertino, and Jorge Lobo. 2017. Shortfall-Based Optimal Placement of Security Resources for Mobile IoT Scenarios. In *Proc. ESORICS*. 419–436.
- [26] Haisheng Tan, Zhenhua Han, Xiang-Yang Li, and Francis C M Lau. 2017. Online Job Dispatching and Scheduling in Edge-Clouds. In *Proc. IEEE INFOCOM*. 1–9.
- [27] Taxi and Limousine Commission (TLC). 2019. 2018 Yellow Taxi Trip Data. <https://data.cityofnewyork.us/Transportation/2018-Yellow-Taxi-Trip-Data/t29m-gskq>
- [28] Liang Tong, Yong Li, and Wei Gao. 2016. A Hierarchical Edge Cloud Architecture for Mobile Computing. In *Proc. IEEE INFOCOM*. 1–9.
- [29] Lin Wang, Lei Jiao, Ting He, Jun Li, and Max Mühlhauser. 2018. Service Entity Placement for Social Virtual Reality Applications in Edge Computing. In *Proc. IEEE INFOCOM*. 468–476.
- [30] Duncan J. Watts and Steven H. Strogatz. 1998. Collective Dynamics of ‘Small-World’ Networks. *Nature* 393, 6684 (jun 1998), 440–442.
- [31] Christian Wolf. 2013. *Advanced Acceleration Techniques for Nested Benders Decomposition in Stochastic Programming*. Ph.D. Dissertation. University of Paderborn.
- [32] Yong Xiao and Marwan Krunz. 2017. QoE and Power Efficiency Tradeoff for Fog Computing Networks with Fog Node Cooperation. In *Proc. IEEE INFOCOM*. 1–9.
- [33] Di Xie, Ning Ding, Y Charlie Hu, and Ramana Kompella. 2012. The Only Constant Is Change: Incorporating Time-Varying Network Reservations in Data Centers. In *Proc. ACM SIGCOMM*. 199–210.
- [34] Jie Xu, Lixing Chen, and Pan Zhou. 2018. Joint Service Caching and Task Offloading for Mobile Edge Computing in Dense Networks. In *Proc. IEEE INFOCOM*. 207–215.
- [35] Yinyu Ye. 1991. An $O(n^3L)$ Potential Reduction Algorithm for Linear Programming. *Math. Program.* 50, 1–3 (mar 1991), 239–258.
- [36] Ruozhou Yu, Guoliang Xue, Vishnu Teja Kilari, and Xiang Zhang. 2018. Deploying Robust Security in Internet of Things. In *Proc. IEEE CNS*. 1–9.
- [37] Ruozhou Yu, Guoliang Xue, and Xiang Zhang. 2018. Application Provisioning in Fog Computing-enabled Internet-of-Things: A Network Perspective. In *Proc. IEEE INFOCOM*. 783–791.
- [38] Deze Zeng, Lin Gu, Song Guo, Zixue Cheng, and Shui Yu. 2016. Joint Optimization of Task Scheduling and Image Placement in Fog Computing Supported Software-Defined Embedded System. *IEEE Trans. Comput.* 65, 12 (dec 2016), 3702–3712.