

Dependent Function Embedding for Distributed Serverless Edge Computing

Shuiguang Deng¹, Senior Member, IEEE, Hailiang Zhao², Zhengzhe Xiang², Member, IEEE, Cheng Zhang, Rong Jiang², Ying Li², Jianwei Yin², Schahram Dustdar³, Fellow, IEEE, and Albert Y. Zomaya⁴, Fellow, IEEE

Abstract—Edge computing is booming as a promising paradigm to extend service provisioning from the centralized cloud to the network edge. Benefit from the development of serverless computing, an edge server can be configured as a carrier of limited serverless functions, in the way of deploying Docker runtime and Kubernetes engine. Meanwhile, an application generally takes the form of directed acyclic graphs (DAGs), where vertices represent dependent functions and edges represent data traffic. The status quo of minimizing the completion time (a.k.a. makespan) of the application motivates the study on optimal function placement. However, current approaches lose sight of proactively splitting and mapping the traffic to the logical data paths between the heterogeneous edge servers, which could affect the makespan significantly. To remedy that, we propose an algorithm, termed as Dependent Function Embedding (DPE), to get the optimal edge server for each function to execute and the moment it starts executing. DPE finds the best segmentation of each data traffic by exquisitely solving several infinity norm minimization problems. DPE is theoretically verified to achieve the global optimality. Extensive experiments on Alibaba cluster trace show that DPE significantly outperforms two baseline algorithms in makespan by 43.19% and 40.71%, respectively.

Index Terms—Edge computing, dependent function embedding, directed acyclic graph, function placement, task scheduling

1 INTRODUCTION

IN recent years, the micro-services application architecture has achieved rapid advances. By changing applications from monolithic to small pieces of code as functions, Function-as-a-Service (FaaS) is leading its way to the future service pattern of cloud computing [1], [2]. Combing FaaS with lightweight containerization and service orchestration tools, such as the Docker runtime [3] and the Kubernetes engine

[4], the concept of serverless computing becomes increasingly popular [5]. Serverless computing offers a platform that allows the execution of software without providing any notion of the underlying computing clusters, operating systems, VMs or containers [6]. It is a one step forward in the abstraction staircase from Infrastructure-as-a-Service (IaaS) to Platform-as-a-Service (PaaS) [2].

Meanwhile, with more and more applications offloaded to remote cloud data-centers, it is hard to meet the QoS requirements of latency-sensitive applications [7]. To mitigate latency, near-data processing within the network edge is a more applicable way to gain insights, which leads to the birth of edge computing. Generally, edge computing refers to leveraging the computation and communication enabled servers, located at the network edge, to make quick response to mobile and IoT applications [8], [9]. Edge servers can be co-located with telecommunication equipment in multiple places across radio access networks (RAN) to the core network (5GC), in the way of small-scale data-centers or machine rooms [10].

Demonstrating common features with the requirements of Internet of Things (IoT) applications, the adaptation of serverless in edge computing has attracted special attention from both industrial community and academia, leading to the birth of serverless edge computing [11]. The paradigm of serverless edge computing allows users to execute their differentiated applications without managing the underlying servers and clusters. Serverless has proved to be more cost-efficient and user-friendly compared with a traditional IaaS architecture in many pilot projects for edge computing [2]. Nevertheless, serverless edge computing faces a series of problems need to be solved urgently. One of the problems restricting its development concerns the scheduling of the functions with

- Shuiguang Deng is with the College of Computer Science and Technology, Zhejiang University, Hangzhou 310058, China, and also with the Institute of Intelligence Applications, Yunnan University of Finance and Economics, Kunming 650221, China. E-mail: dengsg@zju.edu.cn.
- Hailiang Zhao, Cheng Zhang, Ying Li, and Jianwei Yin are with the College of Computer Science and Technology, Zhejiang University, Hangzhou 310058, China. E-mail: {hliangzhao, coolzc, cnliying, zjujw}@zju.edu.cn.
- Zhengzhe Xiang is with Zhejiang University City College, Hangzhou 310015, China. E-mail: xiangzz@zucc.edu.cn.
- Rong Jiang is with the Institute of Intelligence Applications, Yunnan University of Finance and Economics, Kunming 650221, China. E-mail: jiang_rong@aliyun.com.
- Schahram Dustdar is with Distributed Systems Group, Technische Universität Wien, 1040 Vienna, Austria. E-mail: dustdar@dsg.tuwien.ac.at.
- Albert Y. Zomaya is with the School of Computer Science, University of Sydney, Sydney, NSW 2006, Australia. E-mail: albert.zomaya@sydney.edu.au.

Manuscript received 24 Apr. 2021; revised 15 Oct. 2021; accepted 19 Dec. 2021. Date of publication 29 Dec. 2021; date of current version 7 Mar. 2022.

This work was supported in part by the Key Research Project of Zhejiang Province under Grant 2022C01145 and in part by the National Science Foundation of China under Grants U20A20173 and 62125206. Schahram Dustdar's work was supported by the Zhejiang University De-qing Institute of Advanced technology and Industrialization (ZDATI).

(Corresponding author: Ying Li.)

Recommended for acceptance by D. Medhi.

Digital Object Identifier no. 10.1109/TPDS.2021.3137380

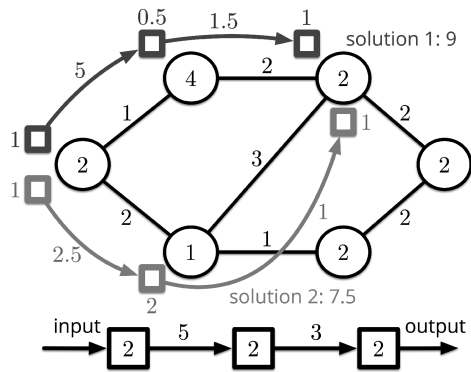


Fig. 1. Two function placement solutions for an SFC with different traffic routing policies.

complex inter-task dependency to the resource-constrained edge [11], [12], [13]. To address this issue, works studying *optimal function placement* across the heterogeneous edge servers are conducted [14], [15], [16]. In these works, applications are structured as a service function chain (SFC) or a directed acyclic graph (DAG) composed of dependent functions, and the placement of each function is obtained by minimizing the makespan of the application, under the trade-off between function processing time and cross-server data transferring overhead.

However, when minimizing the makespan of the application, state-of-the-art approaches only optimize the placement of functions, but how the request to call the function's successors being routed and how the corresponding data stream being mapped onto the virtual links between edge servers are ignored [14], [15]. Actually, routing and management of flow traffic are of great importance for cloud-native applications. In Kubernetes-native systems, Istio is popular for traffic management. It relies on the Envoy proxies co-configured along with the functions [17]. Istio enables the manager of the edge-cloud cluster to configure how each function's request calls its successors, along with its internal output data, routes within an Istio service mesh. Based on flexible and smart configurations, we can find that better utilization of traffic routing and stream mapping can result in less makespan even though the corresponding function placement is not optimal. This phenomenon is captured in Fig. 1. The top half of this figure is an undirected connected graph of six edge servers, abstracted from the physical infrastructure of the heterogeneous edge (In Sections 2 and 3 we will explain how this connected graph is modeled). The numbers tagged in each vertex and beside each edge of the graph are the processing power (measured in *gflop/s*) and guaranteed bandwidth (measured in *GB/s*), respectively. The bottom half is an SFC with three functions. The number tagged inside each function is the required processing power (measured in *gflops*). The number tagged beside each data stream is the size of it (measured in *GB*). Fig. 1 demonstrates two solutions of function placement. The numbers tagged beside vertices and edges of each solution are the time consumed. Just in terms of function placement, solution 1 enjoys lower function processing time ($2.5 < 4$). However, the makespan of solution 2 is 1.5 lesser than solution 1 because it has a better traffic routing policy.

The above example implies that different traffic routing policies could affect the makespan of an application significantly.

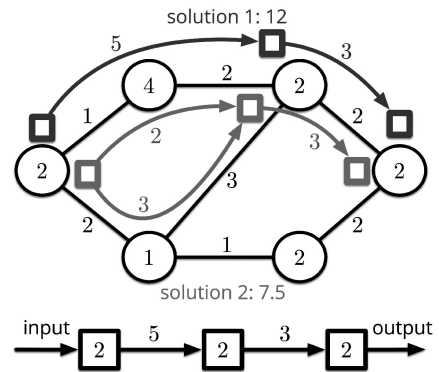


Fig. 2. Two function placement solutions with stream splitting allowed or not, respectively.

It leads us to take traffic routing into consideration *proactively*. In this paper, we name the combination of function placement and stream mapping as *function embedding*. Moreover, if stream splitting is allowed, i.e., the internal output of a function can be split and routed on multiple paths, the makespan decreases further. This phenomenon is captured in Fig. 2. The structure of this figure is the same as Fig. 1. It demonstrates two function embedding solutions with stream splitting allowed or not, respectively. In solution 2, the output stream of the first function is divided into two parts, each with 2 or 3 units. Correspondingly, the times consumed on routing are 3 and 2.5, respectively. Although the two solutions have the same function placement, the makespan of solution 2, which is calculated as $1 + \max\{3, 2.5\} + 1 + 1.5 + 1 = 7.5$, is 4.5 lesser than solution 1.

To capture the importance of proactive traffic routing, in this paper we study the *dependent function embedding* problem with *stream splitting* at the serverless edge. For a DAG with complicated structure, the problem is combinatorial and difficult to solve when the DAG scales up. In this paper, we first present the optimal substructure of the problem. Then, we solve each substructure optimally with dynamic programming. Specifically, for each substructure, we separate several infinity norm minimization sub-problems and solve them optimally by following the analytical solutions. Our paper makes the following contributions:

- *Model contribution:* We study the dependent function embedding problem at the serverless edge. Other than existing works where only function placement is studied, our novel contribution takes proactive traffic routing and data splitting into consideration and leverages dynamic programming as the approach to embed DAGs onto the heterogeneous edge.
- *Algorithm contribution:* We present an algorithm that solves the dependent function embedding problem optimally. We first find the optimal substructure of the problem. In each substructure, we derive the optimal data splitting for the internal data.
- *Experiment contribution:* We conduct extensive experiments on a cluster trace with 2119 unique DAGs from Alibaba [18]. Experimental results show that our algorithm significantly outperforms two algorithms, FixDoc [14] and HEFT [19], on the average completion time by 43.19% and 40.71%, respectively.

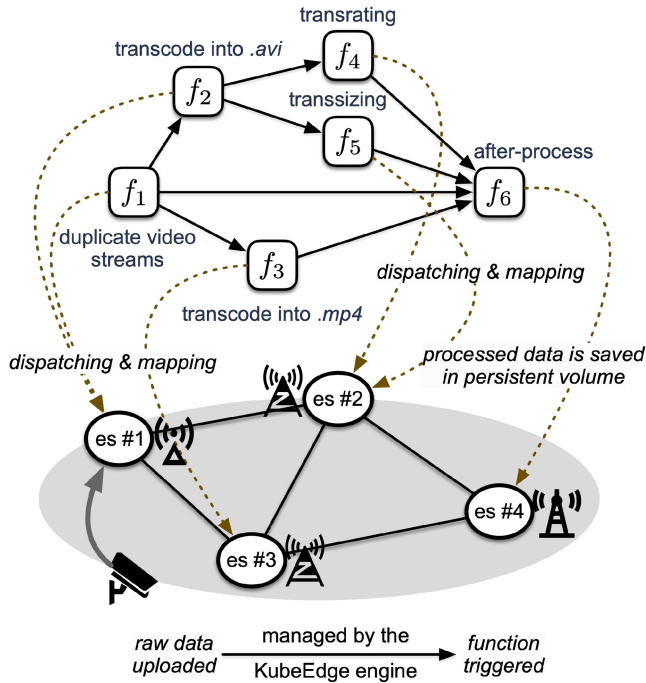


Fig. 3. The architecture of surveillance video processing by leveraging the elastic edge.

The remainder of the paper is organized as follows. In Section 2, we present a working example of the serverless functions. In Section 3, we present system model and formulate the problem. In Section 4, we present the proposed algorithm DPE and several auxiliary algorithms. Performance guarantee and complexity analysis are provided in the same place. The experiment results are demonstrated in Section 5. In Section 6, we review related works on functions placement on the heterogeneous edge. Section 7 concludes this paper.

2 A WORKING EXAMPLE

In this section, we demonstrate a working example on dependent functions at the serverless edge.

The edge network is organized as a *weighted directed graph* [20], [21], where the vertices are edge servers with heterogeneous processing power and the edges are virtual links with certain propagation speed.

The edge network is managed with a lightweight Kubernetes platform, for example, the KubeEdge [22]. Let us deploy an application of surveillance video processing. The procedure is captured in Fig. 3. An edge device, i.e., a surveillance camera, uploads the raw video and pre-prepared configurations¹ to a nearby edge server periodically. With raw video input, the functions are pulled from remote Docker registries and triggered immediately. For video processing, `ffmpeg`² could be used to produce the corresponding Docker images. When it is done, the processed results are saved into a PersistentVolume (PV) configured in the PVC. In the above procedure, the camera only needs to upload the raw video data. The functions will be triggered automatically. The intermediate data produced by

1. A configuration file is used to define Kubernetes object. It should describes the object's name, functions (Docker container instances) to call, PersistentVolumeClaims (PVC), and input & output relations of functions, etc.

2. <https://www.ffmpeg.org/>

TABLE 1
Summary of Key Notations

Notation	Description
$\mathcal{G} \triangleq (\mathcal{N}, \mathcal{L})$	The graph abstracted from the edge network
\mathcal{N}	The set of edge servers
$n_i \in \mathcal{N}$	The i th edge server in \mathcal{G}
$\{\psi_n\}_{n \in \mathcal{N}}$	Processing power of each edge server n
\mathcal{L}	The set of virtual links
$l_{ij} \in \mathcal{L}$	A virtual link from source n_i to target n_j in \mathcal{G}
$\{b_{ij}^{max}\}_{l_{ij} \in \mathcal{L}}$	Maximum bandwidth of the virtual link l_{ij}
$\mathcal{P}(n_i, n_j)$	The set of simple paths from server n_i to n_j
$(\mathcal{F}, \mathcal{E})$	The DAG abstracted from an IoT application
$\{f\}_{f \in \mathcal{F}}$	The set of dependent functions
$c_i, \forall f_i \in \mathcal{F}$	The number of floating point operations of f_i
$e_{ij}, \forall f_i, f_j \in \mathcal{F}$	The data stream from function f_i to f_j
$s_{ij}, \forall f_i, f_j \in \mathcal{F}$	The data stream size from function f_i to f_j
$p(f) \in \mathcal{N}$	The placement of function $f \in \mathcal{F}$
$\varrho \in \mathcal{P}(n_i, n_j)$	A simple path in the set $\mathcal{P}(n_i, n_j)$
z_ϱ	The data stream size route through ϱ
$T(p(f_i))$	The finish time of f_i when scheduled onto $p(f_i)$
$t(e_{ij})$	The time consumed for transferring s_{ij}
$t(p(f_i))$	The time consumed for processing f_i on $p(f_i)$
l_n	The earliest idle time of edge server n
$\sigma(\cdot, \cdot)$	The communication start-up cost
b_{mn}^ϱ	The bandwidth allocated to z_ϱ on l_{mn}

each function is saved into a local volume located at some host path.

To make the most of the quick response of the serverless edge, we need to study where each function is processed and how the flow traffic is mapped, to minimize the completion time as much as possible.

3 SYSTEM MODEL

Let us formulate the heterogeneous edge network as an *undirected connected graph* $\mathcal{G} \triangleq (\mathcal{N}, \mathcal{L})$, where $\mathcal{N} \triangleq \{n_1, \dots, n_N\}$ is the set of edge servers and $\mathcal{L} \triangleq \{l_{ij}\}_{n_i, n_j \in \mathcal{N}}$ is the set of *virtual links*. A virtual link $l_{ij} = (n_i, n_j)$ is an augmented link between edge servers n_i and n_j . In \mathcal{G} , each edge server $n \in \mathcal{N}$ has a processing power ψ_n , measured in *tflop/s* while each virtual link $l_{ij} \in \mathcal{L}$ has a maximum bandwidth b_{ij}^{max} , measured in *GB/s*. We assume $b_{ij}^{max} = b_{ji}^{max}$. When $n_i = n_j$, we simply set the data transferring time as 0 since the intra-server processing is usually negligible. The computation of functions, which is non-preemptive, can be overlapped with communication.

Let us use $\mathcal{P}(n_i, n_j)$ to denote the set of *simple paths*³ from source n_i to target n_j . For arbitrary given edge network, all the simple paths $\{\mathcal{P}(n_i, n_j)\}_{n_i, n_j \in \mathcal{N}}$ can be obtained through depth-first search (DFS). The algorithm designed in this paper needs to know the simple paths between any two edge servers as a priori. In Section 4.4, we give a simple algorithm to calculate them based on DFS.

Key symbols and notations used in this paper are summarized in Table 1.

3.1 Application as a DAG

Each IoT application with dependent functions is modeled as a DAG. Let us use $(\mathcal{F}, \mathcal{E})$ to represent the DAG, where

3. A simple path from n_i to n_j is a path from source n_i to target n_j which contains no loop.

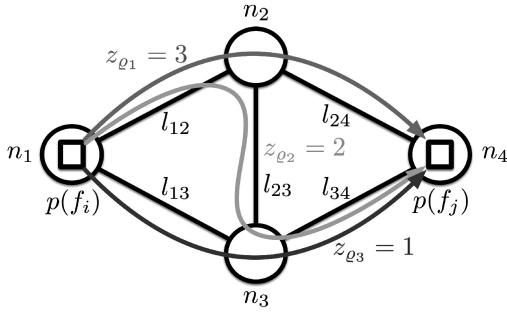


Fig. 4. An example of data stream splitting.

$\mathcal{F} \triangleq \{f_1, \dots, f_F\}$ is the set of F dependent functions listed in topological order⁴. $\forall f_i, f_j \in \mathcal{F}, i \neq j$, if the output stream of f_i is the input of its downstream function f_j , a directed edge e_{ij} exists. $\mathcal{E} \triangleq \{e_{ij} | \forall f_i, f_j \in \mathcal{F}\}$ is the set of all directed edges. For each function $f_i \in \mathcal{F}$, we write c_i for the required number of floating point operations of it. For each directed link $e_{ij} \in \mathcal{E}$, the data stream size is denoted as s_{ij} (measured in GB).

An entry function is a function which does not have predecessors. An exit function is a function which does not have successors. We write \mathcal{F}_{entry} for the set of entry functions and $\mathcal{F}_{exit} \subset \mathcal{F}$ for the set of exit functions of the DAG. In this paper, we make no restrictions on the shape of DAGs. They could be single-entry single-exit, or multi-entry multi-exit.

3.2 Dependent Function Embedding

The dependent function embedding problem is decomposed into two sub-problems, where each function to be dispatched to and how each data stream is mapped onto virtual links.

We write $p(f) \in \mathcal{N}$ for the chosen edge server which $f \in \mathcal{F}$ to be dispatched to. For any function pair (f_i, f_j) and the associated edge $e_{ij} \in \mathcal{E}$, the data stream of size s_{ij} can be splitted and route through different paths in $\mathcal{P}(p(f_i), p(f_j))$. $\forall \varrho \in \mathcal{P}(p(f_i), p(f_j))$, let us use z_ϱ to represent the allocated non-negative data stream size for path ϱ . Then, $\forall e_{ij} \in \mathcal{E}$, we have the following constraint:

$$\sum_{\varrho \in \mathcal{P}(p(f_i), p(f_j))} z_\varrho = s_{ij}.$$

Notice that if $p(f_i) = p(f_j)$, i.e., f_i and f_j are dispatched to the same edge server, then $\mathcal{P}(e_{ij}) = \emptyset$ and the data transferring time is zero.

Fig. 4 gives an example for data splitting. The connected graph has four edge servers and five virtual links. There are four simple paths between n_1 and n_4 . The two squares represent the source function f_i and the destination function f_j . From the edge server $p(f_i)$ to the edge server $p(f_j)$, e_{ij} routes through three out of four simple paths with data size of 3 GB, 2 GB, and 1 GB, respectively. In this example, $s_{ij} = 6$. On closer observation, we can find that two data streams route through l_1 . Each of them is from path ϱ_1 and ϱ_2 with 3 GB and 2 GB, respectively.

4. A topological order of a DAG is a linear ordering of all the vertices such that for every directed edge (f_i, f_j) in this DAG, f_i comes before f_j in this ordering.

Theoretically, overhead exists in the splitting and merging operations of the data stream e_{ij} at the source $p(f_i)$ and the target $p(f_j)$, respectively. However, in our video transcoding scenarios, this overhead is negligible compared with extensive computation of functions or transferring of data streams in gigabytes.

3.3 Involvement Function of Finish Time

Let us use $T(p(f_i))$ to denote the finish time of f_i if it is scheduled to the edge server $p(f_i)$. Considering that the functions of the DAG have dependent relations, for each function $f_j \in \mathcal{F} \setminus \mathcal{F}_{entry}$, $T(p(f_j))$ should involve according to

$$T(p(f_j)) = \max \left\{ \iota_{p(f_j)}, \max_{\forall i: e_{ij} \in \mathcal{E}} \left(T(p(f_i)) + t(e_{ij}) \right) \right\} + t(p(f_j)). \quad (1)$$

In (1), $\iota_{p(f_j)}$ is the earliest idle time of the edge server $p(f_j)$. In our model, the processing of functions is non-preemptive, thus an edge server is idle iff the functions assigned to it complete. $t(e_{ij})$ is the transferring time of data stream e_{ij} and $t(p(f_j))$ is the processing time of f_j on edge server $p(f_j)$. Corresponding to (1), for each entry function $f_i \in \mathcal{F}_{entry}$

$$T(p(f_i)) = \iota_{p(f_i)} + t(p(f_i)), \quad (2)$$

since f_i has no predecessors. In the following, we demonstrate the calculation of $t(e_{ij})$ and $t(p(f_j))$ in turn.

$t(e_{ij})$ is constitutive of two parts, where the first is the communication start-up cost between the two functions f_i and f_j . This cost is mainly decided by the configurations in kube-proxy in this edge-cloud Kubernetes cluster [23]. For example, if we use Envoy to implement the network proxy and communication bus for the cluster, data transfer is mainly handled by Envoy route filters. Before data transfer, Envoy proxy needs to do some preparatory works. For example, looking for the route tables to get the actual routing paths. We simply use $\sigma(f_i, f_j)$ to represent the communication start-up cost. The second is the actual communication cost between $p(f_i)$ and $p(f_j)$, which is actually decided by the slowest data transferring time of z_ϱ among all $\varrho \in \mathcal{P}(p(f_i), p(f_j))$. Considering that the edge network serves for thousands of services and applications, we assume that for each data stream z_ϱ , during its transferring, the bandwidth allocated to it on each virtual link $l_{mn} \in \varrho$ is fixed as b_{mn}^ϱ , and $b_{mn}^\varrho \ll b_{mn}^{max}$ holds. To sum up, $t(e_{ij})$ is defined as

$$t(e_{ij}) \triangleq \sigma(f_i, f_j) + \max_{\varrho \in \mathcal{P}(p(f_i), p(f_j))} \sum_{l_{mn} \in \varrho} \frac{z_\varrho}{b_{mn}^\varrho}. \quad (3)$$

Note that the above formulation ignores the data splitting and merging costs as we have mentioned earlier. In real-world scenarios, the real-time bandwidth for transferring some data stream at some point is unknowable because the network status is always in dynamic changes. Even so, our assumption on the fixed bandwidth allocation is reasonable since it can be guaranteed by the QoS level defined in the generic network slice template (GST) [24]. Besides, our model does not require that each function pair of the DAG enjoys the same bandwidth on each virtual link.

$\forall f_j \in \mathcal{F}$, $t(p(f_j))$ is decided by the processing power of the chosen edge server $p(f_j)$. Since the processing of

functions is non-preemptive, each function is executed with full power. Hence, $t(p(f_j))$ is defined as

$$t(p(f_j)) \triangleq \frac{c_j}{\psi_{p(f_j)}}. \quad (4)$$

3.4 Problem Formulation

After all functions are scheduled, the makespan will be the finish time of the slowest exit function (i.e., the function without successors). Our target is to minimize the makespan of the DAG by finding the optimal $\mathbf{p} \triangleq \{p(f)\}_{\forall f \in \mathcal{F}}$ and the optimal $\mathbf{z} \triangleq \{z_\varrho | \forall \varrho \in \mathcal{P}(p(f_i), p(f_j))\}_{e_{ij} \in \mathcal{E}}$. Thus, the dependent function embedding problem is formulated as

$$\mathbf{P} : \quad \min_{\mathbf{p}, \mathbf{z}} \max_{f \in \mathcal{F}_{exit}} T(p(f))$$

$$s.t. \quad \sum_{\varrho \in \mathcal{P}(p(f_i), p(f_j))} z_\varrho = s_{ij}, \forall e_{ij} \in \mathcal{E}, \quad (5)$$

$$\mathbf{z} \geq \mathbf{0}. \quad (6)$$

4 ALGORITHM DESIGN

In this section, we first give the optimal substructure hidden in \mathbf{P} . Then, we propose the DPE algorithm and provide theoretical analysis on its optimality and complexity. In the end, we give a method to obtain simple paths for the given edge network based on DFS. To simplify the notations, in the following, we replace $\mathcal{P}(p(f_i), p(f_j))$ and $\sigma(p(f_i), p(f_j))$ by \mathcal{P}_{ij} and σ_{ij} , respectively.

4.1 Finding Optimal Substructure

In consideration of the dependency relationship between the fore-and-aft functions, the optimal placement of functions and optimal mapping of data streams cannot be obtained at the same time. Nevertheless, we can solve it optimally step-by-step based on its optimal substructure.

Let us use $T^*(p(f))$ to denote the earliest finish time of function f when it is placed on edge server $p(f)$. Based on (1), $\forall f_j \in \mathcal{F} \setminus \mathcal{F}_{entry}$, we have

$$T^*(p(f_j)) = \max \left\{ \max_{\forall i: e_{ij} \in \mathcal{E}} \left[\min_{p(f_i), \{z_\varrho\}_{\varrho \in \mathcal{P}_{ij}}} \left(T^*(p(f_i)) + t(e_{ij}) \right) \right], t_{p(f_j)} \right\}$$

$$+ t(p(f_j)). \quad (7)$$

Besides, for all the entry functions $f_i \in \mathcal{F}_{entry}$, $T^*(p(f_i))$ is calculated by (2) immediately.

With (7), for each function pair (f_i, f_j) where e_{ij} exists, we can define the sub-problem \mathbf{P}_{sub}

$$\mathbf{P}_{sub} : \min_{p(f_i), \{z_\varrho\}_{\varrho \in \mathcal{P}_{ij}}} \Phi_{ij} \triangleq \left(T^*(p(f_i)) + t(e_{ij}) \right)$$

$$s.t. \quad \sum_{\varrho \in \mathcal{P}_{ij}} z_\varrho = s_{ij}, \quad (8)$$

$$z_\varrho \geq 0, \forall \varrho \in \mathcal{P}_{ij}. \quad (9)$$

Note that (8) and (9) are from a subset of constraints (5) and (6), respectively. In \mathbf{P}_{sub} , we need to decide where f_i is placed and how e_{ij} is mapped. By solving \mathbf{P}_{sub} , we can obtain the earliest finish time $T^*(p(f_j))$ by (7). In this way, \mathbf{P} is solved *optimally* by calculating the earliest finish time of each function in topological order.

4.2 Optimal Data Splitting

To solve \mathbf{P}_{sub} optimally, we first fix the position of f_i , i.e., $p(f_i)$, then concentrate on the optimal mapping of e_{ij} .

To minimize $t(e_{ij})$, let us define a diagonal matrix \mathbf{A} as follows:

$$\mathbf{A} \triangleq \text{diag} \left(\sum_{l_{mn} \in \varrho_1} \frac{1}{b_{mn}^{\varrho_1}}, \sum_{l_{mn} \in \varrho_2} \frac{1}{b_{mn}^{\varrho_2}}, \dots, \sum_{l_{mn} \in \varrho_{|\mathcal{P}_{ij}|}} \frac{1}{b_{mn}^{\varrho_{|\mathcal{P}_{ij}|}}} \right).$$

Obviously, all the diagonal elements of \mathbf{A} are positive real numbers. The variables that need to be determined can be written as $\mathbf{z}_{ij} \triangleq [z_{\varrho_1}, z_{\varrho_2}, \dots, z_{\varrho_{|\mathcal{P}_{ij}|}}]^\top \in \mathbb{R}^{|\mathcal{P}_{ij}|}$. Thus, \mathbf{P}_{sub} is converted into

$$\mathbf{P}_{norm} : \quad \min_{\mathbf{z}_{ij}} \|\mathbf{A}\mathbf{z}_{ij}\|_\infty$$

$$s.t. \quad \begin{cases} \mathbf{1}^\top \mathbf{z}_{ij} = s_{ij}, \\ \mathbf{z}_{ij} \geq \mathbf{0}. \end{cases} \quad (10)$$

We drop $T^*(p(f_i))$ and σ_{ij} readily since constant does not change the optimal solution \mathbf{z}_{ij}^* . \mathbf{P}_{norm} is an infinity norm minimization problem. By introducing slack variables $\tau \in \mathbb{R}$ and $\mathbf{y} \in \mathbb{R}^{|\mathcal{P}_{ij}|}$, \mathbf{P}_{norm} can be transformed into the following *slack form*

$$\mathbf{P}'_{slack} : \quad \min_{\mathbf{z}'_{ij} \triangleq [\mathbf{z}'_{ij}, \mathbf{y}]^\top} \tau$$

$$s.t. \quad \begin{cases} \sum_{\varrho \in \mathcal{P}_{ij}} z_\varrho = s_{ij}, \\ \mathbf{A}\mathbf{z}'_{ij} + \mathbf{y} = \tau \cdot \mathbf{1}, \\ \mathbf{z}'_{ij} \geq \mathbf{0}. \end{cases}$$

\mathbf{P}'_{slack} is feasible and its optimal objective value is finite. As a result, simplex method and interior point method can be applied to obtain the optimal solution efficiently.

However, these standard methods might be unacceptable when the scale of \mathcal{G} increases since simplex method has exponential complexity and interior point method is at least $O(|z'_{ij}|^{3.5})$ in the worst case [25]. Luckily, we can directly obtain the analytical expression of the optimal \mathbf{z}_{ij}^* . The result is introduced in the following theorem.

Theorem 1. *The optimal objective value of \mathbf{P}_{norm} is*

$$\min_{\mathbf{z}_{ij}} \|\mathbf{A}\mathbf{z}_{ij}\|_\infty = \frac{s_{ij}}{\sum_{k=1}^{|\mathcal{P}_{ij}|} 1/A_{k,k}}, \quad (11)$$

iff

$$A_{u,u} z_{ij}^{(u)} = A_{v,v} z_{ij}^{(v)}, 1 \leq u \neq v \leq |\mathcal{P}_{ij}|, \quad (12)$$

where $z_{ij}^{(u)}$ is the u th component of vector \mathbf{z}_{ij} and $A_{u,u}$ is the u th diagonal element of \mathbf{A} .

Proof. For \mathbf{P}_{norm} , we have

$$\|\mathbf{A}\mathbf{z}_{ij}\|_\infty \triangleq \max_k \left\{ |A_{k,k} z_{ij}^{(k)}| \right\} = \lim_{x \rightarrow \infty} \sqrt[x]{\sum_{k=1}^{|\mathcal{P}_{ij}|} \left(A_{k,k} z_{ij}^{(k)} \right)^x},$$

because $\forall k, A_{k,k} > 0, z_{ij}^{(k)} \geq 0$. According to the AM-GM inequality, the following inequality always holds:

$$\frac{\sum_{k=1}^{|\mathcal{P}_{ij}|} \left(A_{k,k} z_{ij}^{(k)} \right)^x}{|\mathcal{P}_{ij}|} \geq \sqrt[|\mathcal{P}_{ij}|]{\prod_{k=1}^{|\mathcal{P}_{ij}|} \left(A_{k,k} z_{ij}^{(k)} \right)^x}, \quad (13)$$

iff (12) is satisfied. It yields that $\forall x > 0$

$$\sqrt[|\mathcal{P}_{ij}|]{\frac{\sum_{k=1}^{|\mathcal{P}_{ij}|} \left(A_{k,k} z_{ij}^{(k)} \right)^x}{|\mathcal{P}_{ij}|}} \geq \sqrt[|\mathcal{P}_{ij}|]{\prod_{k=1}^{|\mathcal{P}_{ij}|} A_{k,k} z_{ij}^{(k)}}. \quad (14)$$

Multiply both sides of (14) by $\sqrt[|\mathcal{P}_{ij}|]{|\mathcal{P}_{ij}|}$, we have

$$\sqrt[|\mathcal{P}_{ij}|]{\sum_{k=1}^{|\mathcal{P}_{ij}|} \left(A_{k,k} z_{ij}^{(k)} \right)^x} \geq \sqrt[|\mathcal{P}_{ij}|]{|\mathcal{P}_{ij}|} \cdot \sqrt[|\mathcal{P}_{ij}|]{\prod_{k=1}^{|\mathcal{P}_{ij}|} A_{k,k} z_{ij}^{(k)}}. \quad (15)$$

By taking the limit of (15), we have

$$\| \mathbf{A} z_{ij} \|_{\infty} \geq \lim_{x \rightarrow \infty} \sqrt[|\mathcal{P}_{ij}|]{|\mathcal{P}_{ij}|} \cdot \sqrt[|\mathcal{P}_{ij}|]{\prod_{k=1}^{|\mathcal{P}_{ij}|} A_{k,k} z_{ij}^{(k)}}. \quad (16)$$

Combining with (10) and (12), the right side of (16) is actually a constant. In other words

$$\begin{aligned} \min_{z_{ij}} \| \mathbf{A} z_{ij} \|_{\infty} &= \lim_{x \rightarrow \infty} \sqrt[|\mathcal{P}_{ij}|]{|\mathcal{P}_{ij}|} \cdot \sqrt[|\mathcal{P}_{ij}|]{\prod_{k=1}^{|\mathcal{P}_{ij}|} A_{k,k} z_{ij}^{(k)}} \\ &= \frac{\sqrt[|\mathcal{P}_{ij}|]{\prod_{k=1}^{|\mathcal{P}_{ij}|} A_{k,k} z_{ij}^{(k)}}}{s_{ij}} \quad \triangleright \text{with (10)} \\ &= \frac{1}{\sum_{k=1}^{|\mathcal{P}_{ij}|} 1/A_{k,k}}. \end{aligned}$$

The result shows that (11) and (12) are the optimal objective value and corresponding optimal condition of \mathbf{P}_{norm} , respectively. \square

Base on (12), we can infer that the optimal variable $z_{ij}^* > \mathbf{0}$ holds, which means that $\forall \varrho \in \mathcal{P}_{ij}$, $z_{\varrho}^* \neq 0$. To sum up, the way to obtain the optimal data splitting and mapping is summarized into Algorithm 1 (OSM) as follows.

Algorithm 1. Optimal Stream Mapping (OSM)

Input: \mathcal{G} , the function pair (f_i, f_j) , and $p(f_j)$
Output: The optimal Φ_{ij}^* , $p^*(f_i)$, and z_{ij}^*

- 1 **for each** $m \in \mathcal{N}$ **do in parallel**
- 2 $p(f_i) \leftarrow m$
- 3 \star Obtain the m th optimal Φ_{ij} by (11) \star /
- 4 $\Phi_{ij}^{(m)} \leftarrow \frac{s_{ij}}{\sum_k 1/A_{k,k}^{(m)}} + T^*(p(f_i))$
- 5 **end for**
- 6 $p^*(f_i) \leftarrow \arg \min_{m \in \mathcal{N}} \Phi_{ij}^{(m)}$
- 7 Calculate z_{ij}^* by (10) and (12) with $\mathbf{A} = \mathbf{A}^{(p^*(f_i))}$

From line 1 to line 5, we can find that OSM solves \mathbf{P}_{norm} by solving $|\mathcal{N}|$ times of \mathbf{P}_{sub} in parallel, each with a different $p(f_i)$. The procedure can be executed in parallel because inter-coupling is nonexistent. In line 4, the objective of \mathbf{P}_{sub} is obtained with the analytical solution (11) directly. The most time-consuming operation lies in line 4 and line 6 since they

have at least one traversal over edge servers and simple paths, respectively. OSM is in $O(\max\{|\mathcal{N}|, |\mathcal{P}_{ij}|\})$ -complexity.

4.3 Dynamic Programming-Based Embedding

Combing OSM with dynamic programming, we have the algorithm DPE (Dynamic Programming-based Embedding).

In DPE, the loop starts from non-entry functions with a topological order. For each non-entry function $f_j \in \mathcal{N} \setminus \mathcal{N}_{entry}$, DPE first fixes its placement $p(f_j)$ as some edge server n in line 3. Then, from line 4 to line 12, DPE solves the sub-problem $\mathbf{P}_{sub}^{(i)}$ by calling OSM for the function pairs (f_i, f_j) , $e_{ij} \in \mathcal{E}$ in turn. If $p^*(f_i)$ has been decided beforehand (under the case where a function is a predecessor of multiple functions), DPE will skip f_i and go to process the next predecessor $f_{i'}$ (line 5 ~ line 7). At the end, DPE updates the finish time of f_j based on the solution of $\{\mathbf{P}_{sub}^{(i)}\}_{\forall e_{ij} \in \mathcal{E}}$ in line 13. Note that the finish time of each entry function f_i should be calculated with (2) before solving $\mathbf{P}_{sub}^{(i)}$. When all the finish time of functions have been calculated, the global minimal makespan of the DAG can be obtained by

$$\max_{f \in \mathcal{F}_{exit}} \arg \min_{p^*} T^*(p^*(f)).$$

The optimal embedding of each function can be retrieved from z^* and p^* .

Algorithm 2. DP-Based Embedding (DPE)

Input: \mathcal{G} and $(\mathcal{F}, \mathcal{E})$
Output: Optimal value and corresponding solution

- 1 **for** $j = |\mathcal{F}_{entry}| + 1$ **to** $|\mathcal{F}|$ **do**
- 2 **for each** $n \in \mathcal{N}$ **do**
- 3 $p(f_j) \leftarrow n$ // Fix the placement of f_j
- 4 **for each** $f_i \in \{f_i | e_{ij} \text{ exists}\}$ **do**
- 5 **if** $p^*(f_i)$ has been decided **then**
- 6 **continue**
- 7 **end if**
- 8 **if** $f_i \in \mathcal{F}_{entry}$ **then**
- 9 $\forall p(f_i) \in \mathcal{N}$, update $T^*(p(f_i))$ by (2)
- 10 **end if**
- 11 Obtain the optimal Φ_{ij}^* , $p^*(f_i)$, and z_{ij}^* by calling OSM
- 12 **end for**
- 13 Update $T^*(p(f_j))$ by (7)
- 14 **end for**
- 15 **end for**
- 16 **return** $\max_{f \in \mathcal{F}_{exit}} \arg \min_{p^*} T^*(p^*(f))$, z^* , and p^*

Fig. 5 demonstrates an example on how PDE works. The top left portion of the figure is a DAG randomly sampled from the Alibaba cluster trace, where all the functions are named in the manner of topological order. $\forall f_i \in \mathcal{F}$, c_i is set as 1. The top right is the edge server cluster \mathcal{G} . The bottom demonstrates how the functions are placed and scheduled by DPE.

In the following, we demonstrate the complexity of DPE.

Theorem 2. *In the worst case, the complexity of DPE is*

$$O(|\mathcal{N}| \times |\mathcal{E}| \times \max\{|\mathcal{N}|, \max_{n_i, n_j} |\mathcal{P}_{ij}|\}).$$

Proof. For each $f_j \in \mathcal{F} \setminus \mathcal{F}_{entry}$, the average number of predecessors is $\frac{|\mathcal{E}|}{|\mathcal{F}_{entry}|}$. Thus, for each placement of f_j , OSM is

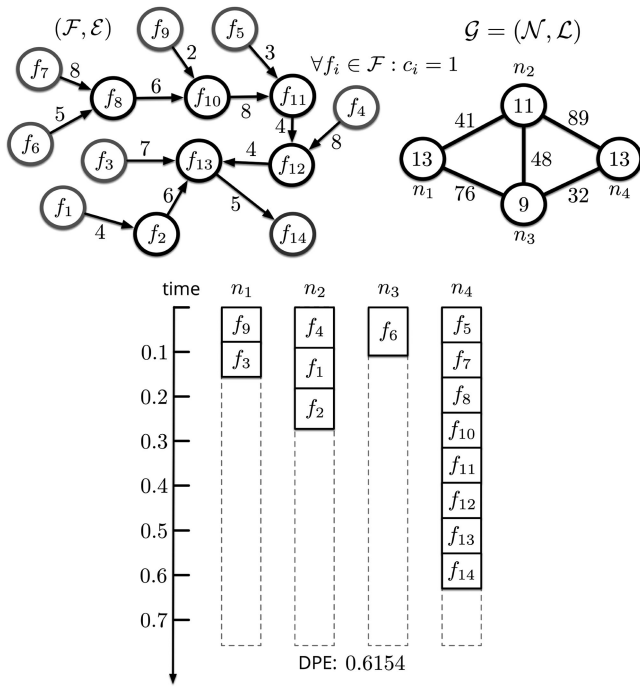


Fig. 5. Embedding of a DAG with the DPE algorithm.

called $\frac{|\mathcal{E}|}{|\mathcal{F} - |\mathcal{F}_{entry}|}$ times on average. Thus, OSM is called

$$\left(|\mathcal{F}| - |\mathcal{F}_{entry}|\right) \times |\mathcal{N}| \times \frac{|\mathcal{E}|}{|\mathcal{F}| - |\mathcal{F}_{entry}|},$$

times in the worst case (under which a function is not a predecessor of more than one function). Due to the complexity of OSM is $O\left(\max\{|\mathcal{N}|, |\mathcal{P}_{ij}|\}\right)$, the result is immediate. \square

4.4 Recursion-Based Path Finding

In this subsection, we demonstrate how to obtain the simple paths for each vertex pair in the undirected graph \mathcal{G} . All the simple paths can be calculated, based on the method proposed in the following, and saved into a shareable table for each edge server. The calculation of simple paths is to be a *once-and-done* operation, and it is calculated immediately after the creation of the Kubernetes cluster. Once the edge network status changes, for example, an edge server is inaccessible or a new edge server is configured into the network, only the related entries in this table need to be re-calculated. The modification has much lower complexity.

Before introducing our algorithm for path finding, we first prove that the simple path finding problem is NP-hard.

Proposition 1. *Finding all the simple paths for all the vertex pairs in the undirected connected graph \mathcal{G} is NP-hard.*⁵

Proof. Suppose there exists a polynomial algorithm \mathcal{A} which can obtain all the simple paths between any source vertex $p(f_i) \in \mathcal{N}$ and any destination vertex $p(f_j) \in \mathcal{N}$. Because \mathcal{G} is a connected graph, $\mathcal{P}(p(f_i), p(f_j))$ is not empty. Since \mathcal{A} lists all the simple paths in $\mathcal{P}(p(f_i), p(f_j))$, we can easily design another simple polynomial

algorithm \mathcal{A}' , based on the results of \mathcal{A} , to judge which simple path in $\mathcal{P}(p(f_i), p(f_j))$ is the longest path. Since \mathcal{G} is undirected and connected, we can claim that this longest simple path has to traverse all vertices of \mathcal{G} . It means that we have found a Hamiltonian Path in polynomial time, which is well known to be NP-hard. As a result, we can conclude that, unless $P = NP$, \mathcal{A} is very unlikely to exist. Therefore, finding all the simple paths between any two vertices of an undirected connected graph is NP-hard. Based on this, we can easily prove that finding all the simple paths for all the vertex pairs of a undirected connected graph is also NP-hard. \square

Let us introduce a new notation, $\Omega(n_i, n_j, \mathcal{M})$, to represent the set of simple paths from n_i to n_j where no path goes through vertices from the set $\mathcal{M} \subseteq \mathcal{N}$. The set of simple paths from n_i to n_j we want, i.e., \mathcal{P}_{ij} , is equal to $\Omega(n_i, n_j, \emptyset)$. $\forall n \in \mathcal{N}$, let us use $\mathcal{A}(n)$ to represent the set of edge servers adjacent to n . Then, $\Omega(n_i, n_j, \mathcal{M})$ can be calculated by the following recursion formula:

$$\Omega(n_i, n_j, \mathcal{M}) = \left\{ J(\varrho, n_i) \mid \bigcup_{m \in \mathcal{S}} \Omega(m, n_j, \mathcal{M} \cup \{n_i\}) \right\},$$

where $\mathcal{S} \triangleq \mathcal{A}(n_i) - \mathcal{M} \cup \{n_i\}$, and $J(\varrho, n_i)$ is a function that joins the node n_i to the path ϱ and returns the new joint path $\varrho \cup \{n_i\}$.

Based on the above recursion formula, we introduce Algorithm 3, RPF (Recursion-based Path Finding). Before calling RPF, we need to initialize the global variables. Specifically, \mathcal{P}_{ij} stores all the simple paths, which is initialized as \emptyset . \mathcal{V} , as the set of visited vertices, is initialized as \emptyset . ϱ , as the path to be calculated currently, which is also initialized as \emptyset .

Algorithm 3. Recursion-Based Path Finding (RPF)

Input: $n_i, n_j \in \mathcal{N}$
Output: \mathcal{P}_{ij}

- 1 $\mathcal{V}, \varrho, \mathcal{P}_{ij} \leftarrow \emptyset, \emptyset, \emptyset$
- 2 $n \leftarrow n_i$
- 3 **if** $n == n_j$ **then**
- 4 $\mathcal{P}_{ij}.append(J(\varrho, n))$ // Store the path $J(\varrho, n)$
- 5 **else**
- 6 $\varrho.push(n); \mathcal{V}.add(n)$
- 7 **for each** $n' \in \mathcal{A}(n) - \mathcal{V}$ **do**
- 8 RPF(n_i, n', n_j) // Recursive call
- 9 **end for**
- 10 $\varrho.pop(); \mathcal{V}.delete(n)$
- 11 **end if**

Although the calling of RPF is a *once-and-done* operation, we still give the complexity of it as follows. Let us use $\kappa(i, j)$ to denote the flops required to compute all the simple paths between n_i and n_j . If \mathcal{G} is fully connected

$$\begin{aligned} \kappa(1, |\mathcal{N}|) &= \sum_{i=2}^{|\mathcal{N}|-1} \left(\kappa(i, |\mathcal{N}|) + 1 \right) + 1 \\ &= (|\mathcal{N}| - 1) + (|\mathcal{N}| - 2) \cdot \kappa(2, |\mathcal{N}|). \end{aligned}$$

5. The proof of this proposition is based on the discussions in the web page <https://stackoverflow.com/questions/9535819/find-all-paths-between-two-graph-nodes>.

Further, we use κ_i to replace $\kappa(i, |\mathcal{N}|)$ by fixing the target vertices as the $|\mathcal{N}|$ th vertex. We can conclude that $\forall i \in \{1, \dots, |\mathcal{N}| - 1\}$

$$\kappa_i = (|\mathcal{N}| - i) + (|\mathcal{N}| - i - 1) \cdot \kappa_{i+1}. \quad (17)$$

Based on the recursion formula (17), we have

$$\begin{aligned} \kappa_1 &= (|\mathcal{N}| - 2)! \cdot \sum_{i=1}^{|\mathcal{N}|-1} \frac{|\mathcal{N}| - i}{(|\mathcal{N}| - i - 1)!} \\ &= (|\mathcal{N}| - 2)! \cdot \sum_{i=1}^{|\mathcal{N}|-1} \frac{(|\mathcal{N}| - i)^2}{(|\mathcal{N}| - i)!} \\ &= (|\mathcal{N}| - 2)! \cdot \sum_{i=1}^{|\mathcal{N}|-1} \frac{i^2}{i!}, \end{aligned} \quad (18)$$

which is the maximum flops required to compute all the simple paths between any two edge servers. To get the upper bound of κ_1 , in the following, we introduce several lemmas.

Lemma 1. $\forall N \geq 7$ and $N \in \mathbb{N}^+$, $N! > N^3(N+1)$.

Proof. The proof is based on *induction*. When $N = 7$, $N! = 5040 > N^3(N+1) = 2744$. The lemma holds. Assume that the lemma holds for $N = q$, i.e., $q! > q^3(q+1)$ (*induction hypothesis*). Then, for $N = q+1$, we have

$$(q+1)! = (q+1) \cdot q! > (q+1)^2 q^3. \quad (19)$$

Notice that the function $g(q) \triangleq (\frac{1}{q} + \frac{2}{q^2} + \frac{4}{q^3})^{-1}$ monotonically increases when $q \in \mathbb{N}^+ - \{1, 2\}$. Hence $g(q) \geq g(3) = \frac{27}{25} > 1$, and

$$1 < \frac{q^3}{(q+2)^2} < \frac{q^3}{(q+1)(q+2)},$$

which indicates that

$$q^3 > (q+1)(q+2). \quad (20)$$

Multiply both sides of (20) by $(q+1)^2$, we get

$$(q+1) \cdot q^3 \cdot (q+1) > (q+1)^3(q+2),$$

which means the induction hypothesis is not violated, and the lemma holds for $q+1$. \square

Lemma 2. $\forall N \geq 2$ and $N \in \mathbb{N}^+$, $\sum_{i=1}^{N-1} \frac{i^2}{i!} < 6 - \frac{1}{N}$.

Proof. We can verify that when $N \in [2, 7] \cap \mathbb{N}^+$, the lemma holds. In the following we prove the lemma holds for $N > 7$ by *induction*. Assume that the lemma holds for $N = q$, i.e., $\sum_{i=1}^{q-1} \frac{i^2}{i!} < 6 - \frac{1}{q}$ (*induction hypothesis*). Then, for $N = q+1$, we have

$$\sum_{i=1}^q \frac{i^2}{i!} < 6 - \frac{1}{q} + \frac{q^2}{q!}. \quad (21)$$

By applying Lemma 1, we get

$$\sum_{i=1}^q \frac{i^2}{i!} < 6 - \frac{1}{q+1},$$

which means the lemma holds for $q+1$. \square

Based on the above lemmas, we can obtain the complexity of RPF, as illustrated in the following theorem:

Theorem 3. *In the worst case, where \mathcal{G} is a fully connected graph and $|\mathcal{N}| \geq 2$, the complexity of RPF is $O((|\mathcal{N}| - 2)!)$.*

Proof. According to Lemma 2

$$\lim_{|\mathcal{N}| \rightarrow \infty} \sum_{i=1}^{|\mathcal{N}|-1} \frac{i^2}{i!} < 6.$$

Hence $\lim_{|\mathcal{N}| \rightarrow \infty} \kappa_1 < 6(|\mathcal{N}| - 2)! = O((|\mathcal{N}| - 2)!)$. \square

In real-world edge computing scenario for IoT stream processing, \mathcal{G} might not be fully connected. Even though, the number of edge servers is small. Thus, the real complexity is much lower. Besides, note that the calling of RPF is actually a *once-and-done* operation. As a priori to DPE, query in this shareable table is of linear complexity. This conclusion is immediate by analyzing line 4 of OSM.

4.5 Extending to Multiple DAGs

DPE can be performed for multiple DAGs. To do this, we only need to concatenate all the DAGs into one augmented DAG. Specifically, if a DAG has more than one entry function, we add a dummy head function f_h and several directed edges $\{e_{hi} | \forall f_i \in \mathcal{F}_{entry}\}$ such that the required floating point operations of f_h is zero and all the s_{hi} are zero. In the same way, if a DAG has more than one exit function, we add a dummy tail function and corresponding directed edges. Based on that, we add a directed edge between the (dummy) tail function of the DAG before and the (dummy) head function of the DAG after. Take the augmented DAG as the input of DPE, we get the embedding results.

5 EXPERIMENTAL VALIDATION

In this section, we conduct extensive experiments to evaluate the effectiveness and efficiency of DPE. Based on a real-world dataset, the Alibaba's cluster trace [18], we first verify the performance of DPE against several popular algorithms on makespan. Then, we analyze the impact of several system parameters.

We summarize the key findings of our experiments as follows, and details can be found in Section 5.2.

- Compared with a well-known heuristic HEFT [19], and a recent algorithm FixDoc [14], DPE achieves the smallest makespan over all the 2119 DAGs with absolute superiority under arbitrary parameters.
- DPE is robust to the system parameters. The advantage of DPE is magnified when the scale of the edge network increases.

5.1 Experiment Setup

IoT Stream Processing Workloads. The simulation is conducted based on Alibaba's cluster trace of data analysis. This dataset contains more than 3 million jobs (called applications in this work), and 20365 jobs with unique DAG information. Considering that there are too many DAGs with only single-digit functions, we sampled 2119 DAGs with different

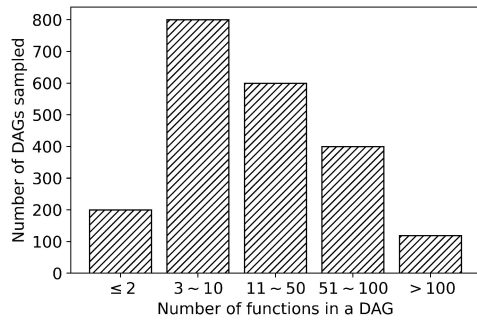


Fig. 6. Data distribution sampled from the cluster trace.

sizes from the dataset. The distribution of the samples is visualized in Fig. 6. For each $f \in \mathcal{F}$, the processing power required and output data size are extracted from the corresponding job in the dataset and scaled to $[1, 10] \times 10^2$ *gflops* and $[5, 15] \times 10^2$ *MB*, respectively.

Heterogeneous Edge Servers. In our simulation, the number of edge servers is 10 in default. Considering that the edge servers are required to formulate a *connected* graph, the impact of the sparsity of the graph is also studied. The processing power of edge servers and the allocated of bandwidth for each data flow are uniformly sampled from $[20, 40]$ *gflops* and $[30, 80] \times 10^2$ *MB/s* in default, respectively.

Algorithms Compared. We compare DPE with the following algorithms.

- *FixDoc* [14]: FixDoc is a function placement and DAG scheduling algorithm with fixed function configuration for the serverless edge computing platform. FixDoc places each function onto *homogeneous* edge servers optimally to minimize the DAG completion time. Actually, [14] also proposes an improved version, GenDoc, with function configuration optimized, too. However, for IoT applications such as video transcoding, on-demand function configuration is not necessary since the tools such as *ffmpeg* can be installed and configured beforehand easily. Thus, we only compare DPE with FixDoc.
- *Heterogeneous Earliest-Finish-Time (HEFT)* [19]: HEFT is a classic heuristic to schedule a set of dependent tasks onto heterogeneous workers with communication time taken into account. Starting with the highest priority, tasks are assigned to different workers to heuristically minimize the overall completion time. HEFT is an algorithm that stands the test of time.

5.2 Experimental Results

All the experiments are implemented in Python 3.7 on macOS Catalina equipped with 3.1 GHz Quad-Core Intel Core i7 and 16 GB RAM. In the following, the unit of the left *y*-axis is *100 seconds*.

5.2.1 Theoretical Performance Verification

Fig. 7 illustrates the overall performance of the three algorithms. For different data batches, DPE can reduce 43.19% and 40.71% of the completion time on average over FixDoc and HEFT on 2119 DAGs. The advantage of DPE is more obvious when the scale of DAG is large because the parallelism is fully

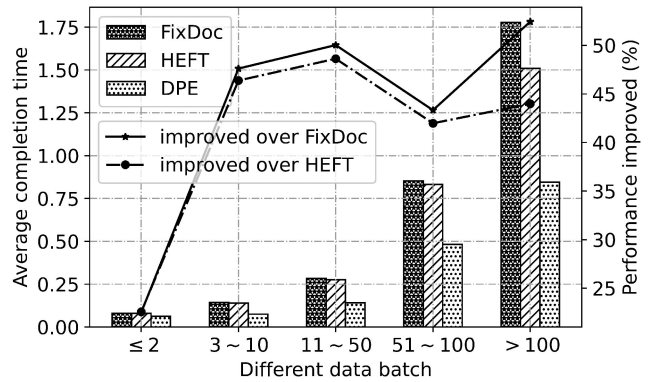


Fig. 7. Average completion time achieved by different algorithms.

guaranteed. Fig. 8 shows the accumulative distribution of 2119 DAGs' completion time. DPE is superior to HEFT and FixDoc on 100% of the DAGs. In Fig. 8, the maximum completion time of a DAG achieved by DPE is 124 seconds. By contrast, only less than 94% of DAGs' completion times achieved by HEFT and FixDoc are smaller than this value.

Figs. 7 and 8 verify the superiority of proactive stream mapping and data splitting. By spreading data streams over multiple virtual links, data transferring time is greatly reduced. Besides, the optimal substructure makes sure DPE can find the optimal placement of each function simultaneously.

5.2.2 Scalability Analysis

Figs. 9 and 10 show the impact of the scale of the heterogeneous edge \mathcal{G} . In Fig. 9, we can find that the average completion time achieved by all algorithms decreases as the edge server increases. The result is obvious because more *idle* servers available, more functions can be executed in parallel without delay. For all data batches, DPE achieves the best result. It is interesting to find that the gap between other algorithms and DPE gets widened when the scale of \mathcal{G} increases. This is because the available simple paths become more and the data transmission time is reduced even further. Fig. 9 also shows the run time of different algorithms at the right *y*-axis. The results show that DPE has the minimum running time overhead.

Fig. 10 shows the impact of sparsity of \mathcal{G} . The horizontal axis is the overall number of simple paths \mathcal{G} . As it increases, \mathcal{G} becomes more denser. Because DPE can reduce transmission

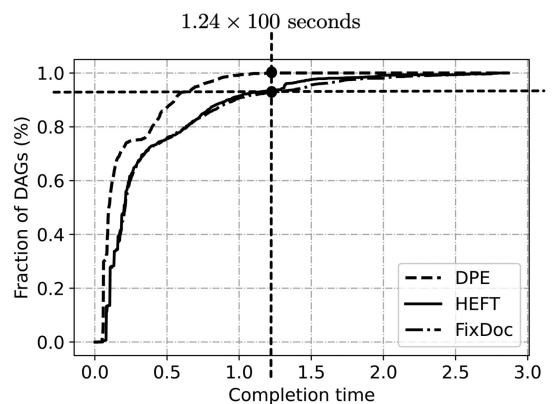


Fig. 8. The accumulative distributions of the completion time achieved by three algorithms, respectively.

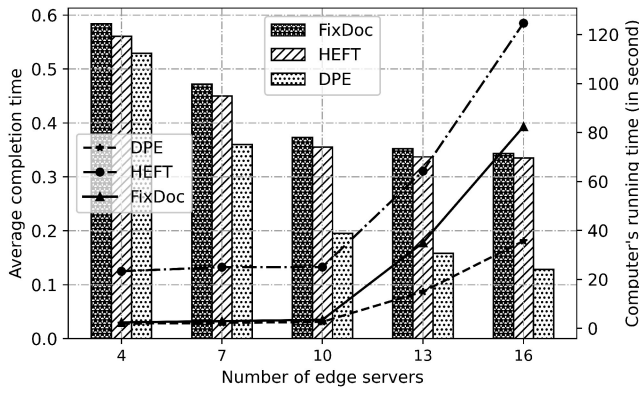


Fig. 9. Average completion time of DAGs and the computer's run time achieved by each algorithm under different numbers of edge servers.

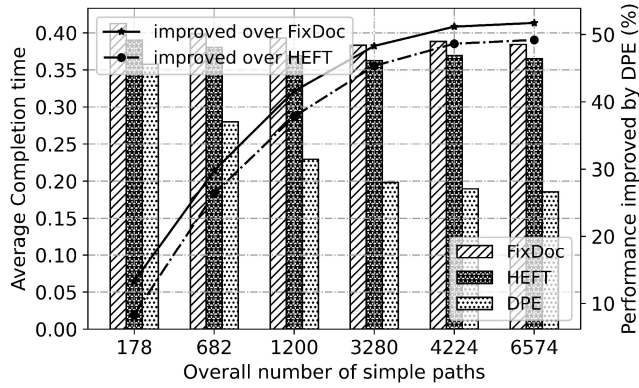


Fig. 10. Average completion time under different sparsities of \mathcal{G} .

time with optimal data splitting and mapping, average completion time achieved by it decreases pretty evident. By contrast, FixDoc and HEFT have no obvious change.

5.2.3 Sensitivity Analysis

Figs. 11 and 12 demonstrate the impact of system parameters, ψ_n and b_l . Notice that $\forall n \in \mathcal{N}, l \in \mathcal{L}, \psi_n$ and b_l are sampled from the interval $[\psi_{lower}, \psi_{upper}]$ and $[b_{lower}, b_{upper}]$ uniformly, respectively. When the processing power and throughput increase, the computation and transmission time achieved by all algorithms are reduced. The results are immediate because a larger processing power directly reduce the computation time of functions while a larger throughput reduces the data transferring time directly. Even so, DPE achieves the smallest average completion time, which verifies the robustness of DPE adequately.

6 RELATED WORKS

In this section, we review related works on function placement and DAG scheduling in serverless edge computing systems.

Studying the optimal function placement is not new. Since cloud computing paradigm became popular, it has been extensively studied in the literature [26], [27], [28]. When bringing function placement into the paradigm of edge computing, especially for the IoT stream processing, different constraints, such as the response time requirement of latency-critical applications, availability of function instances on the heterogeneous edge servers, and the wireless and wired network throughput, etc., should be taken into consideration

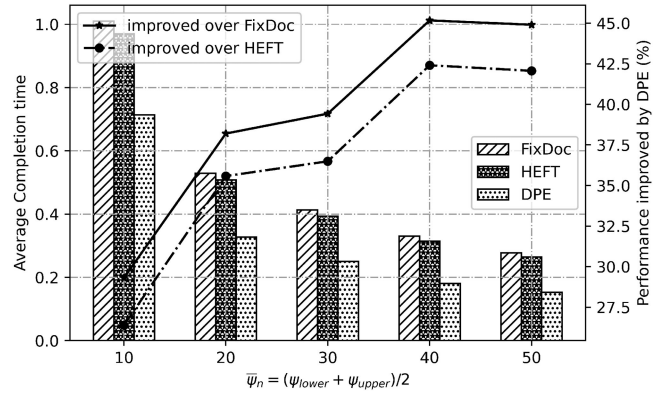


Fig. 11. Average completion time under different processing powers of servers.

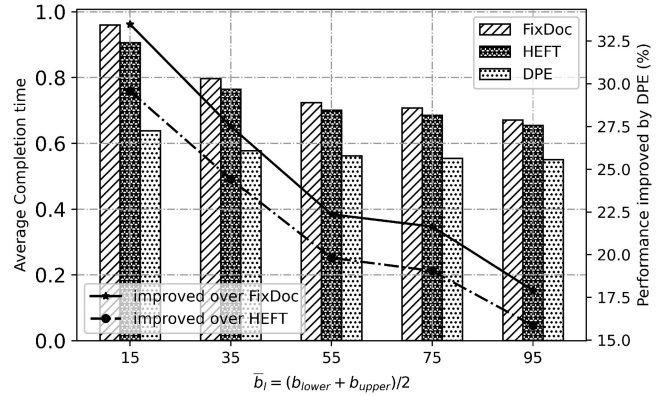


Fig. 12. Average completion time under different throughputs of links.

[29], [30], [31]. In edge computing, the optimal function placement strategy can be used to maximize the network utility [32], minimize the inter-node traffic [33], [34], [35], minimize the makespan of the applications [14], [15], [16], or even minimize the budget of application service providers [36].

In edge computing, the application is either modeled as an individual black-box or a DAG with complicated composite patterns. Considering that the IoT stream processing applications at the edge usually have dependent correlations between the fore-and-aft functions, dependent function placement problem has a strong correlation with DAG dispatching and scheduling. Scheduling algorithms for edge computation tasks have been extensively studied in recent years [19], [37], [38], [39]. In edge computing, the joint optimization of DAG scheduling and function placement is usually NP-hard. As a result, many works can only achieve a near optimal solution based on heuristic or greedy policy. For example, Gedeon *et al.* proposed a heuristic-based solution for function placement across a three-tier edge-fog-cloud heterogeneous infrastructure [40]. Cat *et al.* proposed a greedy algorithm for function placement by estimating the response time of paths in a DAG with queue theory [41]. Although FixDoc [14] can achieve the global optimal function placement, the completion time can be reduced further by optimizing the stream mapping.

7 CONCLUSION

In this paper, we study the optimal dependent function embedding problem at the serverless edge. We first point out

that proactive stream mapping and data splitting could have a strong impact on the makespan of DAGs with several use cases. Based on these observations, we design the DPE algorithm, which is theoretically verified to achieve the global optimality for an arbitrary DAG when the topological order of functions is given. DPE obtains the optimal stream mapping for each function pair with dependent relations. Extensive simulations based on the Alibaba cluster trace dataset verify that our algorithms can reduce the makespan significantly compared with a state-of-the-art function placement and scheduling methods, FixDoc, and a widely accepted algorithm, HEFT. The algorithms proposed in this paper is an off-line algorithm. We leave the extension to online scenarios to our future work.

REFERENCES

- [1] P. Castro, V. Ishakian, V. Muthusamy, and A. Slominski, "The rise of serverless computing," *Commun. ACM*, vol. 62, no. 12, pp. 44–54, Nov. 2019.
- [2] M. S. Aslanpour *et al.*, "Serverless edge computing: Vision and challenges," in *Proc. Australas. Comput. Sci. Week Multiconf.*, 2021, pp. 1–10.
- [3] "Docker: Accelerate how you build, share and run modern applications." Accessed: Oct. 2020. [Online]. Available: <https://www.docker.com/>
- [4] "Kubernetes: Production-grade container orchestration." Accessed: Oct. 2020. [Online]. Available: <https://kubernetes.io/>
- [5] E. Jonas *et al.*, "Cloud programming simplified: A Berkeley view on serverless computing," *CoRR*, vol. abs/1902.03383, 2019.
- [6] E. van Eyk, L. Toader, S. Talluri, L. Versluis, A. Uță, and A. Iosup, "Serverless is more: From PAAS to present cloud computing," *IEEE Internet Comput.*, vol. 22, no. 5, pp. 8–17, Sep./Oct. 2018.
- [7] Y. Mao, C. You, J. Zhang, K. Huang, and K. B. Letaief, "A survey on mobile edge computing: The communication perspective," *IEEE Commun. Surveys Tut.*, vol. 19, no. 4, pp. 2322–2358, Oct.–Dec. 2017.
- [8] P. Porabage, J. Okwuibe, M. Liyanage, M. Ylianttila, and T. Taleb, "Survey on multi-access edge computing for Internet of Things realization," *IEEE Commun. Surveys Tut.*, vol. 20, no. 4, pp. 2961–2991, Oct.–Dec. 2018.
- [9] S. Deng, H. Zhao, W. Fang, J. Yin, S. Dustdar, and A. Y. Zomaya, "Edge intelligence: The confluence of edge computing and artificial intelligence," *IEEE Internet Things J.*, vol. 7, no. 8, pp. 7457–7469, Aug. 2020.
- [10] 5G PPP Architecture Working Group, "View on 5G architecture: Version 3.0," Feb. 2020. [Online]. Available: <https://doi.org/10.5281/zenodo.3265031>
- [11] B. Javadi, J. Sun, and R. Ranjan, "Serverless architecture for edge computing," *Edge Comput.: Models, Technol. Appl.*, pp. 249–264, 2020.
- [12] P. Aditya, I. E. Akkus, A. Beck, R. Chen, V. Hilt, I. Rimac, K. Satzke, and M. Stein, "Will serverless computing revolutionize NFV?" *Proc. IEEE*, vol. 107, no. 4, pp. 667–678, Apr. 2019.
- [13] L. Baresi, D. Figueira Mendonça, and M. Garriga, "Empowering low-latency applications through a serverless edge computing architecture," in *Proc. Serv.-Oriented Cloud Comput.*, F. De Paoli, S. Schulte, and E. Broch Johnsen, Eds. 2017, pp. 196–210.
- [14] L. Liu, H. Tan, S. H.-C. Jiang, Z. Han, X.-Y. Li, and H. Huang, "Dependent task placement and scheduling with function configuration in edge computing," in *Proc. Int. Symp. Qual. Serv.*, 2019, pp. 1–10.
- [15] S. Khare *et al.*, "Linearize, predict and place: Minimizing the makespan for edge-based stream processing of directed acyclic graphs," in *Proc. 4th ACM/IEEE Symp. Edge Comput.*, New York, NY, USA, 2019, pp. 1–14.
- [16] Z. Zhou, Q. Wu, and X. Chen, "Online orchestration of cross-edge service function chaining for cost-efficient edge computing," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 8, pp. 1866–1880, Aug. 2019.
- [17] "Istio: Connect, secure, control, and observe services." Accessed: Oct. 2020. [Online]. Available: <https://istio.io/latest/>
- [18] "Alibaba cluster trace program." Accessed: Oct. 2020. [Online]. Available: <https://github.com/alibaba/clusterdata>
- [19] H. Topcuoglu, S. Hariri, and M.-Y. Wu, "Performance-effective and low-complexity task scheduling for heterogeneous computing," *IEEE Trans. Parallel Distrib. Syst.*, vol. 13, no. 3, pp. 260–274, Mar. 2002.
- [20] X. Foukas, G. Patounas, A. Elmokashfi, and M. K. Marina, "Network slicing in 5G: Survey and challenges," *IEEE Commun. Mag.*, vol. 55, no. 5, pp. 94–100, May 2017.
- [21] S. Vassilaras *et al.*, "The algorithmic aspects of network slicing," *IEEE Commun. Mag.*, vol. 55, no. 8, pp. 112–119, Aug. 2017.
- [22] "Kubeedge: An open platform to enable edge computing." Accessed: Oct. 2020. [Online]. Available: <https://kubeedge.io/en/>
- [23] "Envoy: An open source edge and service proxy, designed for cloud-native applications." Accessed: Oct. 2020. [Online]. Available: <https://www.envoyproxy.io/>
- [24] GSM Association, "Official document ng.116 - generic network slice template v4.0," Nov. 2020. [Online]. Available: <https://www.gsma.com/newsroom/wp-content/uploads//NG.116-v4.0-2.pdf>
- [25] J. Fearnley and R. Savani, "The complexity of the simplex method," in *Proc. 47th Annu. ACM Symp. Theory Comput.*, 2015, pp. 201–208.
- [26] G. T. Lakshmanan, Y. Li, and R. Strom, "Placement strategies for internet-scale data stream systems," *IEEE Internet Comput.*, vol. 12, no. 6, pp. 50–60, Nov./Dec. 2008.
- [27] L. Tom and V. R. Bindu, "Task scheduling algorithms in cloud computing: A survey," in *Inventive Computation Technologies*, S. Smys, R. Bestak, and A. Rocha, Eds. Cham, Switzerland: Springer, 2020, pp. 342–350.
- [28] B. Addis, D. Belabed, M. Bouet, and S. Secci, "Virtual network functions placement and routing optimization," in *Proc. IEEE 4th Int. Conf. Cloud Netw.*, 2015, pp. 171–177.
- [29] H. Badri, T. Bahreini, D. Grosu, and K. Yang, "Energy-aware application placement in mobile edge computing: A stochastic optimization approach," *IEEE Trans. Parallel Distrib. Syst.*, vol. 31, no. 04, pp. 909–922, Apr. 2020.
- [30] M. Nardelli, V. Cardellini, V. Grassi, and F. Presti, "Efficient operator placement for distributed data stream processing applications," *IEEE Trans. Parallel Distrib. Syst.*, vol. 30, no. 08, pp. 1753–1767, Aug. 2019.
- [31] Z. Ning *et al.*, "Distributed and dynamic service placement in pervasive edge computing networks," *IEEE Trans. Parallel Distrib. Syst.*, vol. 32, no. 06, pp. 1277–1292, Jun. 2021.
- [32] M. Leconte, G. S. Paschos, P. Mertikopoulos, and U. C. Kozat, "A resource allocation framework for network slicing," in *Proc. INFO-COM Conf. Comput. Commun.*, 2018, pp. 2177–2185.
- [33] J. Xu, Z. Chen, J. Tang, and S. Su, "T-storm: Traffic-aware online scheduling in storm," in *Proc. IEEE 34th Int. Conf. Distrib. Comput. Syst.*, 2014, pp. 535–544.
- [34] L. Liu, S. Guo, G. Liu, and Y. Yang, "Joint dynamical VNF placement and SFC routing in NFV-enabled SDNS," *IEEE Trans. Netw. Serv. Manage.*, vol. 18, no. 4, pp. 4263–4276, Dec. 2021.
- [35] S. Yang, F. Li, S. Trajanovski, X. Chen, Y. Wang, and X. Fu, "Delay-aware virtual network function placement and routing in edge clouds," *IEEE Trans. Mobile Comput.*, vol. 20, no. 2, pp. 445–459, Feb. 2021.
- [36] L. Chen, J. Xu, S. Ren, and P. Zhou, "Spatio temporal edge service placement: A bandit learning approach," *IEEE Trans. Wireless Commun.*, vol. 17, no. 12, pp. 8388–8401, Dec. 2018.
- [37] Y. Kao, B. Krishnamachari, M. Ra, and F. Bai, "Hermes: Latency optimal task assignment for resource-constrained mobile computing," *IEEE Trans. Mobile Comput.*, vol. 16, no. 11, pp. 3056–3069, Nov. 2017.
- [38] S. Sundar and B. Liang, "Offloading dependent tasks with communication delay and deadline constraint," in *Proc. IEEE INFO-COM Conf. Comput. Commun.*, 2018, pp. 37–45.
- [39] J. Meng, H. Tan, C. Xu, W. Cao, L. Liu, and B. Li, "Dedas: Online task dispatching and scheduling with bandwidth constraint in edge computing," in *Proc. IEEE Conf. Comput. Commun.*, 2019, pp. 2287–2295.
- [40] J. Gedeon, M. Stein, L. Wang, and M. Muehlhaeuser, "On scalable in-network operator placement for edge computing," in *Proc. 27th Int. Conf. Comput. Commun. Netw.*, 2018, pp. 1–9.
- [41] X. Cai, H. Kuang, H. Hu, W. Song, and J. Lü, "Response time aware operator placement for complex event processing in edge computing," in *Service-Oriented Computing*, C. Pahl, M. Vukovic, J. Yin, and Q. Yu, Eds. Cham, Switzerland: Springer, 2018, pp. 264–278.



Shuiguang Deng (Senior Member, IEEE) received the BS and PhD degree in computer science from the College of Computer Science and Technology, Zhejiang University, China, in 2002 and 2007, respectively. He is currently a full professor with the College of Computer Science and Technology, Zhejiang University. He was visiting scholar with the Massachusetts Institute of Technology in 2014 and Stanford University in 2015. He has authored or coauthored more than 100 papers in journals and refereed conferences. His research interests

include edge computing, service computing, cloud computing, and business process management. He is currently an associate editor for the journal *IEEE Transactions on Services Computing*, *Knowledge and Information Systems*, *Computing* and *IET Cyber-Physical Systems: Theory & Applications*. He was the recipient of Rising Star Award by IEEE TCSVC in 2018. He is a fellow of IET.



Hailiang Zhao received the BS degree from the School of computer science and technology, Wuhan University of Technology, Wuhan, China, in 2019. He is currently working toward the PhD degree with the College of Computer Science and Technology, Zhejiang University, Hangzhou, China. His research interests include edge computing, service computing, and machine learning. He was the recipient of Best Student Paper Award of IEEE ICWS 2019.



Zhengzhe Xiang (Member, IEEE) received the BS and PhD degrees in computer science and technology from Zhejiang University, Hangzhou, China. He was a visiting student with Karlstad University, Sweden, in 2018. He is currently a lecturer with Zhejiang University City College, Hangzhou, China. His research interests include service computing, cloud computing, and edge computing.



Cheng Zhang received the MS degree in 2013 in electrical engineering from Zhejiang University, China, where he is currently working toward the PhD degree in computer science and technology. His research interests include edge computing and edge intelligence.



Rong Jiang received the PhD degree in system analysis and integration from the School of Software, Yunnan University, China. He is currently the deputy dean with the Institute of Intelligence Applications, a distinguished professor, and the doctoral supervisor with the Yunnan University of Finance and Economics, Kunming, China. He is an excellent professional and technical talent with outstanding contributions in Yunnan Province, an expert enjoying special government allowances in Yunnan Province, a young and middle-aged academic and

technical leader and an excellent teacher in Yunnan Province. He is the director of Key Laboratory of Service Computing and security management of Yunnan Provincial Universities, and the director of Kunming Key Laboratory of Information Economy and Information Management. He has authored or coauthored more than 60 papers. His research interests include cloud computing, big data, block chain, AI application and information management, digital economy, and software engineering. He was the recipient of more than 70 prizes in recent years.



Ying Li is currently an associate professor with the College of Computer Science and Technology, Zhejiang University, Hangzhou, China. His research interests include service computing, cloud computing, and data science.



Jianwei Yin received the PhD degree in computer science from Zhejiang University (ZJU) in 2001. He was a visiting scholar with the Georgia Institute of Technology. He is currently a full professor with the College of Computer Science, ZJU. He has authored or coauthored more than 100 papers in top international journals and conferences. His research interests include service computing and business process management. He is an associate editor for *IEEE Transactions on Services Computing*.



Schahram Dustdar (Fellow, IEEE) is currently a full professor of computer science informatics with a focus on internet technologies heading the Distributed Systems Group, TU Wien. Since December 2016, he has been the chairman of the Informatics Section of the Academia Europaea. From 2004 to 2010, he was an honorary professor of information systems with the Department of Computing Science, University of Groningen (RuG), The Netherlands. From December 2016 until January 2017, he was a visiting professor with the University of Sevilla, Spain, and from January until June 2017, he was a visiting professor with UC Berkeley, USA. He has been a member of the IEEE Conference Activities Committee (CAC) since 2016, of the Section Committee of Informatics of the Academia Europaea since 2015, and a member of the Academia Europaea: The Academy of Europe, Informatics Section since 2013. He was the recipient of the ACM Distinguished Scientist Award (2009) and the IBM Faculty Award (2012). He is an associate editor for *IEEE Transactions on Services Computing*, *ACM Transactions on the Web*, and *ACM Transactions on Internet Technology*. He is on the editorial board of *IEEE Internet Computing*. He is the editor-in-chief of *Computing* (an SCI-ranked journal of Springer).



Albert Y. Zomaya (Fellow, IEEE) is currently the chair professor of high-performance computing and networking with the School of Computer Science and the director of the Centre for Distributed and High-Performance Computing with the University of Sydney. He has authored or coauthored more than 600 scientific papers and articles and is the co-author or the editor of more than 30 books. As a sought-after speaker, he has delivered more than 190 keynote addresses, invited seminars, and media briefings.

His research interests include parallel and distributed computing and complex systems. He is currently the editor in chief of the *ACM Computing Surveys* and was the editor in chief of *IEEE Transactions on Computers* from 2010 to 2014 and *IEEE Transactions on Sustainable Computing* from 2016 to 2020. He is a decorated scholar with numerous accolades including the fellowship of the American Association for the Advancement of Science, and the Institution of Engineering and Technology, U.K. He is an elected fellow of the Royal Society of New South Wales and an elected foreign member of Academia Europaea. He was the recipient of 1997 Edgeworth David Medal from the Royal Society of New South Wales for outstanding contributions to Australian Science, the IEEE Technical Committee on Parallel Processing Outstanding Service Award (2011), IEEE Technical Committee on Scalable Computing Medal for Excellence in Scalable Computing (2011), IEEE Computer Society Technical Achievement Award (2014), ACM MSWIM Reginald A. Fessenden Award (2017), and the New South Wales Premier's Prize of Excellence in Engineering and Information and Communications Technology (2019).

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.