

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/302919660>

# Optimal Workload Allocation in Fog-Cloud Computing Towards Balanced Delay and Power Consumption

Article in IEEE Internet of Things Journal · January 2016

DOI: 10.1109/JIOT.2016.2565516

CITATIONS

422

READS

3,604

5 authors, including:



**Ruilong Deng**

Zhejiang University

88 PUBLICATIONS 3,070 CITATIONS

[SEE PROFILE](#)



**Rongxing Lu**

University of New Brunswick

366 PUBLICATIONS 14,652 CITATIONS

[SEE PROFILE](#)



**Chengzhe Lai**

Xi'an University of Posts and Telecommunications

26 PUBLICATIONS 1,172 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



biosensors [View project](#)



Heckman-Opdam theory and multivariable hypergeometric functions [View project](#)

# Optimal Workload Allocation in Fog-Cloud Computing Towards Balanced Delay and Power Consumption

Ruilong Deng, *Member, IEEE*, Rongxing Lu, *Senior Member, IEEE*, Chengzhe Lai, *Member, IEEE*, Tom H. Luan, *Member, IEEE*, and Hao Liang, *Member, IEEE*

**Abstract**—Mobile users typically have high demand on localized and location-based information services. To always retrieve the localized data from the remote cloud, however, tends to be inefficient, which motivates fog computing. The fog computing, also known as edge computing, extends cloud computing by deploying localized computing facilities at the premise of users, which pre-stores cloud data and distributes to mobile users with fast-rate local connections. As such, fog computing introduces an intermediate fog layer between mobile users and cloud, and complements cloud computing towards low-latency high-rate services to mobile users. In this fundamental framework, it is important to study the interplay and cooperation between the edge (fog) and the core (cloud). In this paper, the tradeoff between power consumption and transmission delay in the fog-cloud computing system is investigated. We formulate a workload allocation problem which suggests the optimal workload allocations between fog and cloud towards the minimal power consumption with the constrained service delay. The problem is then tackled using an approximate approach by decomposing the primal problem into three subproblems of corresponding subsystems, which can be respectively solved. Finally, based on simulations and numerical results, we show that by sacrificing modest computation resources to save communication bandwidth and reduce transmission latency, fog computing can significantly improve the performance of cloud computing.

**Index Terms**—Cloud computing, fog computing, optimization, power consumption-delay tradeoff, workload allocation.

## I. INTRODUCTION

Manuscript received January 17, 2016; accepted May 04, 2016. This work was supported in part by EEE Cybersecurity Research Program at Nanyang Technological University, Alberta Innovates Technology Futures (AITF) post-doctoral fellowship, International Science and Technology Cooperation and Exchange Plan in Shaanxi Province, China (2015KW-010), National Natural Science Foundation of China Research Grant 61502386, and a research grant from the Natural Science and Engineering Research Council (NSERC) of Canada. R. Lu would like to thank the support from Nanyang Technological University's College of Engineering Proposal Preparatory Grant and MOE Tier 1 (M4011450). A preliminary version was presented at IEEE ICC 2015 [1]. The review of this paper was coordinated by Prof. Andrea Zanella. Paper no. IoT-0846-2016. (*Corresponding author: Rongxing Lu.*)

R. Deng is with the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore 639798; he is now also with the Department of Electrical and Computer Engineering, University of Alberta, Edmonton, AB, Canada T6G 1H9 (e-mail: ruilong@ualberta.ca).

R. Lu is with the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore 639798 (e-mail: rxlu@ntu.edu.sg).

C. Lai is with the National Engineering Laboratory for Wireless Security, Xi'an University of Posts and Telecommunications, Xi'an, 710121, China (e-mail: lc.zu@xupt.com).

T. H. Luan is with the School of Information Technology, Deakin University, Burwood, Victoria 3125, Australia (e-mail: tom.luan@deakin.edu.au).

H. Liang is with the Department of Electrical and Computer Engineering, University of Alberta, Edmonton, AB, Canada T6G 1H9 (e-mail: hao2@ualberta.ca).

Copyright (c) 2012 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

THE INTERNET has shifted to the cloud based structure. As reported in Cisco Cloud Index (2013-2018), since 2008, most Internet traffic has originated or terminated in a data center. By 2016, it is predicted that nearly two-thirds of total workloads in traditional IT space will be processed in the cloud. However, with the surging mobile traffic generated in recent years, the transmission of the extraordinarily huge-volume data to the cloud has not only posed a heavy burden on communication bandwidth, but also resulted in unbearable transmission latency and degraded service to end users [2]–[4]. In addition to real-time interaction and low latency, with mobile users and traffic becoming dominant nowadays, the support of mobility and geo-distribution is also critical [5]–[7]. Therefore, with cloud becoming the overarching approach for centralized information storage, retrieval, and management, and mobile devices becoming the major destination of information, the successful integration of cloud computing and mobile applications therefore represents an important task.

To address the above challenges, Cisco has delivered the concept of fog computing in 2014, which aims to process in part workload and services locally on fog devices (such as hardened routers, switches, IP video cameras, etc.), rather than being transmitted to the cloud [8]. This is by introducing a new intermediate fog layer between mobile users and cloud as shown in Fig. 1. The fog layer is composed of geo-distributed fog servers which are deployed at the edge of networks, e.g., parks, bus terminals, shopping centers, etc. Each fog server is a highly virtualized computing system, similar to a light-weight cloud server, and is equipped with the on-board large-volume data storage, compute, and wireless communication facility. The fog servers bridges the mobile users and cloud. On one hand, fog servers directly communicate with the mobile users through single-hop wireless connections using the off-the-shelf wireless interfaces, such as WiFi, Bluetooth, etc. With the on-board compute facility and pre-cached contents, they can independently provide pre-defined service applications to mobile users without assistances from cloud or Internet. On the other hand, the fog servers can be connected to the cloud so as to leverage the rich functions and application tools of the cloud. Therefore, “the fog is a cloud close to the ground”. Fog computing is not to substitute but to complement cloud computing, in order to ease bandwidth burden and reduce transmission latency. In particular, the fog can support and facilitate applications that do not fit well with the cloud: (i) applications that require very low and predictable latency, such as online gaming and video conferencing; (ii) geographically distributed applications such as pipeline monitoring and sensor networks; (iii) fast mobile applications such as smart con-

nected vehicles; and (iv) large-scale distributed control systems such as smart energy distribution and smart traffic lights [9]–[12].

While the fog provides localization, i.e., enabling the real-time interaction and low latency at the network edge, the cloud provides centralization, the integration of which inspires applications that require the interplay and cooperation between the edge (fog) and the core (cloud), particularly for big data and Internet of Things [13]–[16]. From this perspective, we showcase some specific use cases of fog-cloud computing [17], [18]. For example, fog devices deployed inside a multi-floor shopping center can deliver delay-sensitive services including indoor navigation and flyers distribution to mobile users through WiFi, and forward delay-tolerant requests such as feedback statistical analysis to cloud servers for centralized processing. Fog devices deployed at a park lot can provide the pre-cached information including park maps and local accommodations, and, by connecting to cloud servers, send timely alerts and notifications to drivers. Fog devices deployed inside an inter-state bus can deliver onboard video streaming and social networking services to passengers using WiFi. The onboard fog devices connect to cloud servers through cellular networks to refresh the pre-cached contents and update application services, and also report users' data such as their feedbacks to cloud servers for centralized processing.

In this paper, we consider a fog-cloud computing system. On one hand, with the huge-volume and ever-increasing service requests, the power consumption on powering up (and cooling) cloud servers is soaring. It is thus important and desirable to consider the energy management in the fog-cloud computing system [19], [20]. On the other hand, it is equally crucial to guarantee the quality of service (e.g., latency requirements) of end users. The reason is that the unbearable response latency leads to revenue loss of service providers since end users will subscribe to other vendors with better service [21]. To this end, we systematically investigate the fundamental tradeoff between the power consumption and delay in the fog-cloud computing system.

In this paper, firstly, we model the power consumption function and delay function of each part of the fog-cloud computing system, and formulate the workload allocation problem. Then, we develop an approximate approach to solve the primal problem through decomposition, and formulate three subproblems of three corresponding subsystems. These subproblems can be respectively solved via existing optimization techniques. Finally, based on simulations and numerical results, we show that fog computing can significantly improve the performance of cloud computing in terms of reducing communication latency. To the best of our knowledge, this is an early effort towards providing a systematic framework of computation and communication co-design in the fog-cloud computing system. We hope that this pioneering work can throw light on how the fog can extend and complement the cloud. Specifically, the original contributions of this paper are summarized in the following three folds:

- 1) We cast a mathematical framework to investigate the power consumption-delay tradeoff problem by workload allocation in the fog-cloud computing system.
- 2) We develop an approximate approach to decompose the primal problem into three subproblems of corresponding subsystems, and solve them respectively.
- 3) We conduct extensive simulations to demonstrate that the fog can significantly complement the cloud with much reduced communication latency.

Beyond the low latency characteristic as addressed in this paper, the possible advantages of a fog architecture include mobility support, geo-distribution, and location/context awareness [22], [23]. Not only can the geo-distributed fog device infer its own location, but also the fog device can track end users' devices to support mobility, which would be a game changing factor for location-based services and applications. Besides, the geo-distribution can also provide rich network context information, such as the local network condition, traffic statistics, and client status information, which can be used by fog applications to offer context-aware optimization.

The remainder of this paper is organized as follows. The related works are introduced in Section II. We describe the model of the fog-cloud computing system and formulate the power consumption-delay tradeoff problem in Section III. In Section IV, we approximately decompose the primal problem into three subproblems of corresponding subsystems. Simulations are conducted in Section V with numerical results, and concluding remarks are drawn in Section VI with future work.

## II. RELATED WORKS

Cloud computing, a kind of Internet-based paradigm, refers to both the applications delivered as services over the Internet and the hardware and software in the data centers that provide these services [24], [25]. The research on cloud computing has attracted great attention with a large quantity of literatures. For example, Armbrust *et al.* [26] quantify comparisons between cloud and conventional computing, and identify the top technical and non-technical obstacles and opportunities of cloud computing. The emergence of cloud computing has established a trend towards building massive, energy-hungry, and geographically distributed Internet data centers as cloud servers. Due to their enormous energy consumption, Rao *et al.* [19], [21] investigate how to coordinate the collection of data centers so as to minimize the electricity expense while maintaining the quality of the cloud computing service. Our work extends from the existing related papers on cloud computing to a newly emerged paradigm named fog computing. However, the transition is not trivial, since fog is quite different from cloud in terms of location, distribution, and computing capability.

On the other hand, fog computing, characterized by extending cloud computing to the network edge, has become a buzzword today [22], [23]. With similar frameworks such as cloudlet, follow me cloud, and edge computing, fog computing receives considerable attention recently. For example, Bonomi *et al.* [10] define the characteristics of fog computing which make it an appropriate platform for a number of critical services and applications in Internet of Things and big data analytics. Stojmenovic *et al.* [27], [28] review a handful of literatures that expand the applications of fog computing

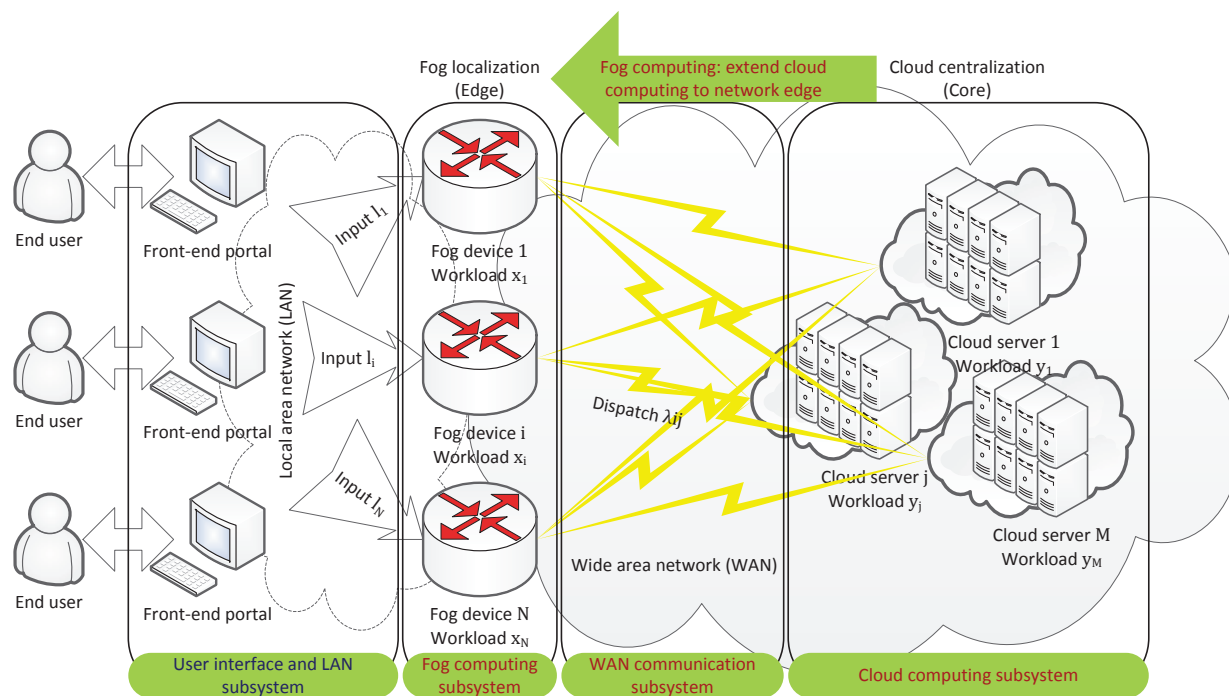


Fig. 1. An overall architecture of a fog-cloud computing system with four subsystems and their interconnections/interactions.

to a series of real scenarios, such as smart grid, vehicular networks, cyber-physical systems, etc. Security and privacy issues are further disclosed according to current fog computing paradigm. Since the fog is not to substitute but to complement the cloud, it is worthy of studying the interaction and cooperation between them. However, existing methodologies need to be changed to accommodate the bi-layer fog-cloud model. To our knowledge, a systematic framework of computation and communication co-design does not seem to be studied so far in the context of fog-cloud. Our work serves as a starting point to address this issue, in which we study the tradeoff between power consumption and delay in the fog-cloud computing system.

### III. SYSTEM MODEL AND PROBLEM FORMULATION

We illustrate an overall architecture of the fog-cloud computing system in Fig. 1, which has been divided into four subsystems. The front-end portals act as user interfaces that receive service requests from end users. These requests are separately input to a set  $\mathcal{N}$  of fog devices through a local area network (LAN). Since fog devices are generally located in the vicinity of end users, thus the LAN communication delay could be omitted (compared to WAN). Fog computing can process some of the delay-sensitive requests and forward others to cloud computing [29]. There is a set  $\mathcal{M}$  of cloud servers, each of which hosts a number of homogeneous computing machines. The unprocessed requests are dispatched from each fog device to each cloud server through a wide area network (WAN). Since WAN covers a large geographical area from the edge throughout to the core, the communication delay and constrained bandwidth should be taken into account. In the following, we mainly consider the power consumption

and computation/communication delay of the latter three subsystems (i.e., fog computing, WAN communication, and cloud computing). Some important notations used in this paper are summarized in Table I. In the rest of this work, we also use the following mathematical notations from linear algebra:  $\mathbf{x}^T$  denotes the transpose of  $\mathbf{x}$ ;  $\mathbf{1}$  denotes the all-ones vector; and  $\mathbf{0}$  denotes the all-zeros vector.

TABLE I  
SUMMARY OF NOTATIONS

Symbol	Definition	Unit <sup>a</sup>
$i, N, \mathcal{N}$	index, number, set of fog devices	n/a
$j, M, \mathcal{M}$	index, number, set of cloud servers	n/a
$l_i$	traffic arrival rate to fog device $i$	#(requests)/s
$x_i$	workload assigned to fog device $i$	#(requests)/s
$\lambda_{ij}$	traffic rate dispatched from fog device $i$ to cloud server $j$	#(requests)/s
$y_j$	workload assigned to cloud server $j$	#(requests)/s
$L$	total input from all front-end portals	#(requests)/s
$X$	workload allocated for fog computing	#(requests)/s
$Y$	workload allocated for cloud computing	#(requests)/s
$P$	power consumption	unit power
$D$	delay	unit time
$\bar{D}$	system delay constraint	unit time
$v_i$	service rate at fog device $i$	#(requests)/s
$f_j$	machine CPU frequency at cloud server $j$	#(cycles)/s
$\sigma_j$	binary: on/off state of cloud server $j$	n/a
$n_j$	integer: machine number at cloud server $j$	n/a
$d_{ij}$	communication delay from fog device $i$ to cloud server $j$	unit time
$\eta_i$	weighting factor at fog device $i$	n/a
$\bar{D}_j$	delay threshold at cloud server $j$	unit time

<sup>a</sup>The unit of a quantity may be omitted in the rest of the paper if it is specified here.



## A. System Model

1) *Power Consumption of Fog Device*: For the fog device  $i$ , the computation power consumption can be modelled by a function of the computation amount  $x_i$ , which is a monotonic increasing and strictly convex function. The piece-wise linear function and quadratic function are two alternatives [30]. In fact, the fog computing devices can accommodate any form of power consumption functions as long as they satisfy the following two properties: (i) the computation power consumption always increases as the computation amount increases; (ii) the marginal power consumption for each fog device is increasing. For simplicity but without loss of generality, we can express the power consumption  $P_i^{\text{fog}}$  of the fog device  $i$  by the following function of the computation amount  $x_i$ :

$$P_i^{\text{fog}} \triangleq a_i x_i^2 + b_i x_i + c_i,$$

where  $a_i > 0$  and  $b_i, c_i \geq 0$  are pre-determined parameters.

2) *Computation Delay of Fog Device*: Assuming a queueing system, for the fog device  $i$  with the traffic arrival rate  $x_i$  and service rate  $v_i$ , the computation delay (waiting time plus service time)  $D_i^{\text{fog}}$  is

$$D_i^{\text{fog}} \triangleq \frac{1}{v_i - x_i}.$$

3) *Power Consumption of Cloud Server*: Each cloud server hosts a number of homogeneous computing machines. The configurations (e.g., CPU frequency) are assumed to be equal for all machines at the same server. Thus, each machine at the same server has the same power consumption profile. We approximate the power consumption value of each machine at the cloud server  $j$  by a function of the machine CPU frequency  $f_j$ :  $A_j f_j^p + B_j$ , where  $A_j$  and  $B_j$  are positive constants, and  $p$  varies from 2.5 to 3 [21].

When the allocated workload increases, more cloud servers are powered on; while when it decreases, the excess servers are turned off for energy saving [31]. Let a binary variable  $\sigma_j$  denote the on/off state of the cloud server  $j$ , where 1 means that the server is on and 0 means off. Besides, let an integer variable  $n_j$  denote the number of turned-on machines at the cloud server  $j$ . Thus, the power consumption  $P_j^{\text{cloud}}$  of the cloud server  $j$  can be obtained by multiplying the on/off state, the on-state machine number, and each machine power consumption value [19]:

$$P_j^{\text{cloud}} \triangleq \sigma_j n_j (A_j f_j^p + B_j).$$

4) *Computation Delay of Cloud Server*: The M/M/n queueing (or Erlang-C) model is employed to characterize each cloud server. In this model, the computation delay (waiting time plus service time) is  $\left[ \frac{C(n, \lambda/\mu)}{n\mu - \lambda} + \frac{1}{\mu} \right]$ , where  $n$  is the number of machines,  $\lambda$  and  $\mu$  are the traffic arrival rate and service rate respectively, and  $C(n, \lambda/\mu)$  is the Erlang's C formula [32, Ch. 2]. At the cloud server  $j$ , assume that each machine has the same service rate  $\mu_j$ . We can generally convert  $\mu_j$  to  $f_j$  by  $\mu_j = f_j/K$ , where  $K$  is in terms of #cycles/request.

From the above, for the cloud server  $j$  with the on/off state  $\sigma_j$  and  $n_j$  turned-on machines, when each machine has the

traffic arrival rate  $y_j$  and service rate  $f_j/K$  respectively, the computation delay  $D_j^{\text{cloud}}$  is given by

$$D_j^{\text{cloud}} \triangleq \sigma_j \left[ \frac{C(n_j, y_j K / f_j)}{n_j f_j / K - y_j} + \frac{K}{f_j} \right].$$

5) *Communication Delay for Dispatch*: Let  $d_{ij}$  denote the delay of the WAN transmission path from the fog device  $i$  to the cloud server  $j$ . Thus, when the traffic rate dispatched from the fog device  $i$  to the cloud server  $j$  is  $\lambda_{ij}$ , the corresponding communication delay  $D_{ij}^{\text{comm}}$  is

$$D_{ij}^{\text{comm}} \triangleq d_{ij} \lambda_{ij}.$$

## B. Constraints

1) *Workload Balance Constraint*: Let  $L$  denote the total request input from all front-end portals. The traffic arrival rate from all front-end portals to the fog device  $i$  is denoted by  $l_i$ . Thus, we have

$$L \triangleq \sum_{i \in \mathcal{N}} l_i.$$

Besides, let  $X$  and  $Y$  denote the workload allocated for fog computing and cloud computing, respectively. Then, we have

$$\begin{cases} X \triangleq \sum_{i \in \mathcal{N}} x_i \\ Y \triangleq \sum_{j \in \mathcal{M}} y_j. \end{cases}$$

We describe the workload balance constraint on the traffic rate dispatched from each fog device to each cloud server. The end-user requests are either handled by a fog device, or forwarded to a cloud server to be processed. The corresponding relationships between the workload and traffic rate are listed as (i) workload balance constraint for each fog device:

$$l_i - x_i = \sum_{j \in \mathcal{M}} \lambda_{ij} \quad \forall i \in \mathcal{N}, \quad (1)$$

(ii) workload balance constraint for each cloud server:

$$\sum_{i \in \mathcal{N}} \lambda_{ij} = y_j \quad \forall j \in \mathcal{M}. \quad (2)$$

From (i) and (ii) we can easily obtain (iii) workload balance constraint for the holistic fog-cloud computing system:

$$L = X + Y.$$

2) *Fog Device Constraint*: For the fog device  $i$ , there exists a limit on the processing ability due to physical constraints. Let  $x_i^{\text{max}}$  denote the computation capacity of the fog device  $i$ . In addition, the workload  $x_i$  assigned to the fog device  $i$  should be no more than the traffic arrival rate  $l_i$  to that device. From the above, we have

$$0 \leq x_i \leq \min \{x_i^{\text{max}}, l_i\} \quad \forall i \in \mathcal{N}. \quad (3)$$

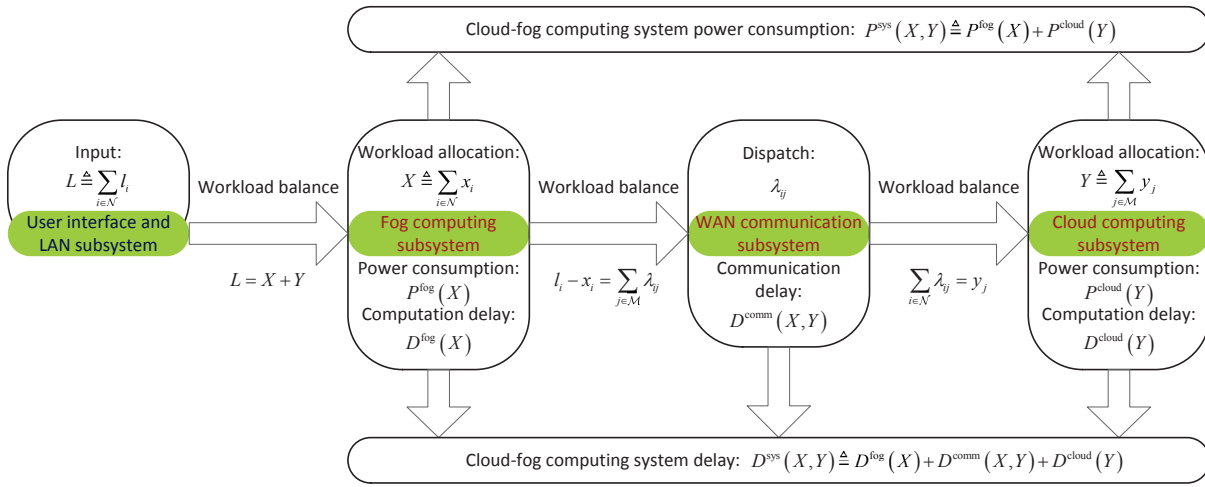


Fig. 2. An overall framework of power consumption-delay tradeoff by workload allocation in a fog-cloud computing system.

3) *Cloud Server Constraint:* For the cloud server  $j$ , firstly, we have

$$y_j \geq 0 \quad \forall j \in \mathcal{M}. \quad (4)$$

Besides, there exists a limit on the computation rate of each machine due to physical constraints. Let  $f_j^{\min}$  and  $f_j^{\max}$  denote the lower and upper bound on the machine CPU frequency, respectively:

$$f_j^{\min} \leq f_j \leq f_j^{\max} \quad \forall j \in \mathcal{M}. \quad (5)$$

In addition, for the cloud server  $j$ , the number of machines  $n_j$  has an upper bound  $n_j^{\max}$ . Thus, for the integer variable  $n_j$ , we have

$$n_j \in \{0, 1, 2, \dots, n_j^{\max}\} \quad \forall j \in \mathcal{M}. \quad (6)$$

Finally, the binary variable  $\sigma_j$  denote the on/off state of the cloud server  $j$ . When  $\sigma_j$  equals 1, it means that the cloud server  $j$  is on; when  $\sigma_j$  equals 0, it means that the cloud server  $j$  is off, and meanwhile the number of on-state machines equals 0. Thus, we have

$$\sigma_j \in \{0, 1\} \quad \forall j \in \mathcal{M}. \quad (7)$$

4) *WAN Communication Bandwidth Constraint:* For simplicity but without loss of generality, the traffic rate  $\lambda_{ij}$  is assumed to be dispatched from the fog device  $i$  to the cloud server  $j$  through one transmission path. Furthermore, these transmission paths do not overlap with each other. There is a limitation  $\lambda_{ij}^{\max}$  on the bandwidth capacity of each path. Thus, the bandwidth constraint of the WAN communication is

$$0 \leq \lambda_{ij} \leq \lambda_{ij}^{\max} \quad \forall i \in \mathcal{N}, \forall j \in \mathcal{M}. \quad (8)$$

### C. Problem Formulation

Towards the power consumption-delay tradeoff in fog-cloud computing, on one hand, it is important and desirable to minimize the aggregated power consumption of all fog devices and cloud servers. The power consumption function of the fog-cloud computing system is defined as

$$P^{\text{sys}} \triangleq \sum_{i \in \mathcal{N}} P_i^{\text{fog}} + \sum_{j \in \mathcal{M}} P_j^{\text{cloud}}.$$

On the other hand, it is equally crucial to guarantee the quality of service (e.g., latency requirements) of end users. The end-user experienced delay consists of the computation (including queueing) delay and communication delay. Therefore, the delay function of the fog-cloud computing system is defined as

$$D^{\text{sys}} \triangleq \sum_{i \in \mathcal{N}} D_i^{\text{fog}} + \sum_{j \in \mathcal{M}} D_j^{\text{cloud}} + \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} D_{ij}^{\text{comm}}.$$

We consider the problem of minimizing the power consumption of the fog-cloud computing system while guaranteeing the required delay constraint  $\bar{D}$  for end users. That is, we have the Primal Problem (PP):

$$\begin{aligned} \min_{x_i, y_j, \lambda_{ij}, f_j, n_j, \sigma_j} \quad & P^{\text{sys}} \\ \text{s.t.} \quad & \begin{cases} D^{\text{sys}} \leq \bar{D} \\ (1) - (8). \end{cases} \end{aligned}$$

The decision variables are the workload  $x_i$  assigned to the fog device  $i$ , the workload  $y_j$  assigned to the cloud server  $j$ , the traffic rate  $\lambda_{ij}$  dispatched from the fog device  $i$  to the cloud server  $j$ , as well as the machine CPU frequency  $f_j$ , the machine number  $n_j$ , and the on/off state  $\sigma_j$  at the cloud server  $j$ . The objective of workload allocation in the fog-cloud computing system is to tradeoff between (i) the system power consumption and (ii) the end-user experienced delay.

## IV. DECOMPOSITION AND SOLUTION

Note that in PP, the decision variables come from different subsystems and are tightly coupled with each other, which makes the relationship between the workload allocation and the power consumption-delay tradeoff not clear. To address this issue, we develop an approximate approach to decompose PP into three subproblems of corresponding subsystems, which can be respectively solved via existing optimization techniques. We illustrate the decomposition and each subproblem/subsystem interactions in Fig. 2, which provides an overall framework of power consumption-delay tradeoff by workload allocation in the fog-cloud computing system.

### A. Power Consumption-Delay Tradeoff for Fog Computing

We consider to tradeoff between the power consumption and computation delay in the fog computing subsystem. That is, we have the Subproblem One (**SP1**):

$$\begin{aligned} \min_{x_i} \quad & \sum_{i \in \mathcal{N}} \left( a_i x_i^2 + b_i x_i + c_i + \frac{\eta_i}{v_i - x_i} \right) \\ \text{s.t.} \quad & \begin{cases} \sum_{i \in \mathcal{N}} x_i = X \\ (3), \end{cases} \end{aligned}$$

where the adjustable parameter  $\eta_i$  is a weighting factor to tradeoff between the power consumption and computation delay at the fog device  $i$ .

Given the workload  $X$  allocated for the fog computing subsystem, **SP1** is a convex problem with linear constraints. This problem can be easily solved using convex optimization techniques such as interior-point methods [33]–[35, Ch. 11]. After we obtain the optimal workload  $x_i^*$  assigned to the fog device  $i$ , we can calculate the power consumption and computation delay in the fog computing subsystem respectively as

$$\begin{cases} P^{\text{fog}}(X) = \sum_{i \in \mathcal{N}} [a_i (x_i^*)^2 + b_i x_i^* + c_i] \\ D^{\text{fog}}(X) = \sum_{i \in \mathcal{N}} \frac{1}{v_i - x_i^*}. \end{cases}$$

### B. Power Consumption-Delay Tradeoff for Cloud Computing

At the cloud server  $j$ , for the delay-sensitive requests, their response delay should be bounded by a certain threshold that is specified as the service level agreement, since the agreement violation would result in loss of business revenue. We assume that the response delay should be smaller than an adjustable parameter  $\bar{D}_j$ , which can be regarded as the delay threshold that identifies the revenue/penalty region at the cloud server  $j$ :

$$D_j^{\text{cloud}} \leq \bar{D}_j.$$

We consider to tradeoff between the power consumption and computation delay in the cloud computing subsystem. That is, we have the Subproblem Two (**SP2**):

$$\begin{aligned} \min_{y_j, f_j^*, n_j^*, \sigma_j} \quad & \sum_{j \in \mathcal{M}} \sigma_j n_j (A_j f_j^p + B_j) \\ \text{s.t.} \quad & \begin{cases} \sum_{j \in \mathcal{M}} y_j = Y \\ D_j^{\text{cloud}} \leq \bar{D}_j \quad \forall j \in \mathcal{M} \\ (4) - (7). \end{cases} \end{aligned}$$

Given the workload  $Y$  allocated for the cloud computing subsystem, **SP2** is a mixed integer nonlinear programming (MINLP) problem, which is generally difficult to tackle. Since the generalized Benders decomposition (GBD) is an effective method to solve this problem with guaranteed optimality, we design the GBD algorithm in Appendix A [36]–[38, Ch. 13]. After we obtain the optimal workload  $y_j^*$  assigned to the cloud server  $j$  and the optimal solution  $f_j^*$ ,  $n_j^*$ , and  $\sigma_j^*$ , we can calculate the power consumption and computation delay in

the cloud computing subsystem respectively as

$$\begin{cases} P^{\text{cloud}}(Y) = \sum_{j \in \mathcal{M}} \sigma_j^* n_j^* [A_j (f_j^*)^p + B_j] \\ D^{\text{cloud}}(Y) = \sum_{j \in \mathcal{M}} D_j^{\text{cloud}*} = \sum_{j \in \mathcal{M}} \sigma_j^* \bar{D}_j. \end{cases}$$

### C. Communication Delay Minimization for Dispatch

We consider the traffic dispatch rate  $\lambda_{ij}$  to minimize the communication delay in the WAN subsystem. That is, we have the Subproblem Three (**SP3**):

$$\begin{aligned} \min_{\lambda_{ij}} \quad & \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} d_{ij} \lambda_{ij} \\ \text{s.t.} \quad & (1), (2), \text{ and } (8). \end{aligned}$$

From Section IV-A and IV-B, given the workload  $X$  allocated for fog computing and  $Y$  for cloud computing, we can obtain the optimal workload  $x_i^*$  assigned to the fog device  $i$  and  $y_j^*$  assigned to the cloud server  $j$ . Given  $x_i^*$  and  $y_j^*$ , **SP3** is regarded as an assignment problem. Since this problem can be efficiently solved using the Hungarian method in polynomial time, we design the Hungarian algorithm in Appendix B [39]. After we obtain the optimal traffic rate  $\lambda_{ij}^*$  dispatched from the fog device  $i$  to the cloud server  $j$ , we can calculate the communication delay in the WAN subsystem as

$$D^{\text{comm}}(X, Y) = \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} d_{ij} \lambda_{ij}^*.$$

### D. Putting It All Together

Based on the above decomposition and the solution to three subproblems, on one hand, the power consumption function of the fog-cloud computing system is rewritten as

$$P^{\text{sys}}(X, Y) \triangleq P^{\text{fog}}(X) + P^{\text{cloud}}(Y),$$

which means that the system power consumption comes from the fog devices and cloud servers. On the other hand, the delay function of the fog-cloud computing system is rewritten as

$$D^{\text{sys}}(X, Y) \triangleq D^{\text{fog}}(X) + D^{\text{cloud}}(Y) + D^{\text{comm}}(X, Y),$$

which means that the system delay comes from the computation delay of the fog devices and cloud servers, as well as the communication delay of the WAN.

After solving the above three subproblems, we can approximately solve **PP** by considering the following approximate problem named **PP-approx**:

$$\begin{aligned} \min_{X, Y} \quad & P^{\text{sys}}(X, Y) \\ \text{s.t.} \quad & \begin{cases} D^{\text{sys}}(X, Y) \leq \bar{D} \\ X + Y = L, \end{cases} \end{aligned}$$

which can be iteratively solved. The approximation ratio is dependent on the choice of two adjustable parameters  $\eta_i$  and  $\bar{D}_j$ . If these parameters could be chosen appropriately, then the solution to **PP-approx** would be the optimal solution to **PP**. How to evaluate the approximation ratio of the proposed decomposition is left as our future work.

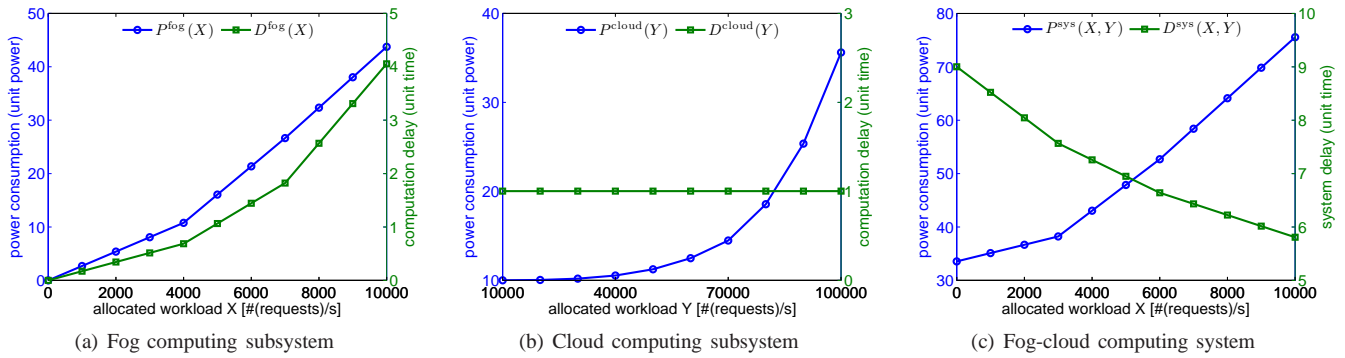


Fig. 3. An illustration of power consumption-delay tradeoff by workload allocation in a fog-cloud computing system.

## V. NUMERICAL RESULTS

Simulation results are presented in this section to validate the power consumption-delay tradeoff by workload allocation to fog computing and cloud computing. For simplicity but without loss of generality, we consider the scenario with five fog devices and three cloud servers (Internet data centers) in the fog-cloud computing system. It can be extended to more fog devices and more cloud servers, with the similar results. Some important parameters used in the simulation are summarized in Table II, referring to [19], [21]. The following results are obtained by MATLAB.

TABLE II  
PARAMETER SETUP

Parameter	Value	Parameter	Value
$l_i$	$[3 \ 1.5 \ 1.5 \ 2 \ 2] \times 10^4$	$f_j^{\min}$	1.0
$A_j$	$[3.206 \ 4.485 \ 2.370]$	$f_j^{\max}$	$[3.4 \ 2.4 \ 3.0]$
$B_j$	$[68 \ 53 \ 70]$	$n_j^{\max}$	$[3 \ 6 \ 2.5] \times 10^4$
$p, K$	3, 1	$\frac{1}{D_j}$	unit time

Firstly, we vary the workload  $X$  allocated for fog computing from 0 to  $10^4$ , to evaluate how they affect the power consumption  $P^{\text{fog}}(X)$  and computation delay  $D^{\text{fog}}(X)$  in the subsystem. Under different values of  $X$ , we solve **SP1** and obtain the optimal workload  $x_i^*$  assigned to the fog device  $i$ . Based on this we calculate  $P^{\text{fog}}(X)$  and  $D^{\text{fog}}(X)$ , and draw their curves in Fig. 3(a). It is seen that both the power consumption and computation delay increase with the workload allocated for fog computing.

Then, we vary the workload  $Y$  allocated for cloud computing from  $10^4$  to  $10^5$ , to evaluate how they affect the power consumption  $P^{\text{cloud}}(Y)$  and computation delay  $D^{\text{cloud}}(Y)$  in the subsystem. Under different values of  $Y$ , we solve **SP2** and obtain the optimal workload  $y_j^*$  assigned to the cloud server  $j$ . Based on this we calculate  $P^{\text{cloud}}(Y)$  and  $D^{\text{cloud}}(Y)$ , and draw their curves in Fig. 3(b). The result shows that the computation delay stays steady while the power consumption increases with the workload allocated for cloud computing.

Finally, based on the above  $x_i^*$  and  $y_j^*$ , we further solve **SP3** and obtain the communication delay  $D^{\text{comm}}(X, Y)$  in the WAN subsystem. Based on these we calculate the system power consumption  $P^{\text{sys}}(X, Y)$  and delay  $D^{\text{sys}}(X, Y)$ , and

draw their curves in Fig. 3(c). From the numerical results, we note that the power consumption of fog devices dominates the system power consumption, while the communication delay of the WAN dominates the system delay. Therefore, when the fog workload is low, the fog power consumption is low and so is the system power consumption, while the WAN communication delay is high and so is the system delay, and vice versa. The figure illustrates that, when some of workload is allocated for fog computing, the system delay decreases while the system power consumption increases. This is because in the fog-cloud computing system, cloud computing is more powerful and energy-efficient than fog computing; while the fog, with the advantage of physical proximity to end users, can sacrifice modest computation resources to save WAN bandwidth and reduce communication latency, in such a way to significantly improve the performance of the cloud.

## VI. CONCLUSION

In this paper, we have introduced the vision of fog computing, a newly emerged paradigm that extends cloud computing to the edge of the network. Concretely, we develop a systematic framework to investigate the power consumption-delay tradeoff issue in the fog-cloud computing system. We formulate the workload allocation problem and approximately decompose the primal problem into three subproblems, which can be respectively solved within corresponding subsystems. Simulation and numerical results are presented to show the fog's complement to the cloud. We hope that this pioneering work can provide guidance on studying the interaction and cooperation between the fog and cloud.

Note that in this paper the optimization is performed in a centralized manner. For the future work, we intend to further consider the case that the optimization is performed in a distributed manner. In that case, the required information exchange and communication overhead need to be carefully investigated.

## APPENDIX A

### SOLVE **SP2** USING GBD ALGORITHM

**Definition 1:** define  $\mathbf{y}, \mathbf{f}, \mathbf{n}, \boldsymbol{\sigma}$  as the vectors of  $y_j, f_j, n_j, \sigma_j$ , and  $Y, F, N, \Sigma$  as the definition domains of  $y_j, f_j, n_j, \sigma_j$ , i.e., (4)-(7).



We now follow [38, Ch. 13] to solve **SP2** using the GBD algorithm. For MINLP **SP2**,  $\mathbf{y}$  and  $\mathbf{f}$  are continuous, while  $\mathbf{n}$  and  $\boldsymbol{\sigma}$  are integer variables. Let  $\mathbf{y}^*$ ,  $\mathbf{f}^*$ ,  $\mathbf{n}^*$ , and  $\boldsymbol{\sigma}^*$  denote the optimal solution. Clearly, finding the optimal integer variables  $\mathbf{n}^*$  and  $\boldsymbol{\sigma}^*$  is the critical part of solving MINLP. When the integer variables are determined, MINLP reduces to a linear programming (LP) problem, which is generally easy to tackle. In other words, once  $\mathbf{n}^*$  and  $\boldsymbol{\sigma}^*$  are determined,  $\mathbf{y}^*$  and  $\mathbf{f}^*$  can be easily solved.

The GBD algorithm is an iterative approach for solving MINLP and the underlying intuition is described as follows. MINLP is decomposed into a master problem (MP) and a subproblem (SP). MP is an integer programming problem, which aims to determine the integer variables by considering only the integer constraints (with lower bound solution  $LB$ ). When the integer variables are determined, MINLP reduces to SP (an LP problem) to determine the continuous variables (with upper bound solution  $UB$ ). In general cases the determined integer variables are not optimal, but they can be improved by adding new integer constraints into MP, such that the search/feasible space shrinks and the newly determined integer variables gradually approach the optimum. For example, then SP has the feasible solution but  $UB > LB$  (i.e.,  $\mathbf{n}$  and  $\boldsymbol{\sigma}$  are not optimal), in order to improve  $\mathbf{n}$  and  $\boldsymbol{\sigma}$ , the new  $LB$  should be large than previous  $LB$ s, by adding the feasibility constraint (9a) into MP. When SP is infeasible to solve, in order to avoid obtaining the improper  $\mathbf{n}$  and  $\boldsymbol{\sigma}$  again, the infeasibility constraint (9b) is added into MP. The optimal solution is converged when  $|UB - LB| \leq \epsilon$ , where  $\epsilon$  is error tolerance (stopping criterion). The iterative approach is summarized in **Algorithm 1**, which involves the following definitions.

**Definition 2:** objective function  $\mathcal{F}(\mathbf{f}, \mathbf{n}, \boldsymbol{\sigma})$  and constraint functions  $\mathcal{G}(\mathbf{y})$ ,  $\mathcal{H}(\mathbf{y}, \mathbf{f}, \mathbf{n}, \boldsymbol{\sigma})$ :

$$\begin{cases} \mathcal{F}(\mathbf{f}, \mathbf{n}, \boldsymbol{\sigma}) \triangleq \sum_{j \in \mathcal{M}} \sigma_j n_j (A_j f_j^p + B_j) \\ \mathcal{G}(\mathbf{y}) \triangleq \sum_{j \in \mathcal{M}} y_j - Y \\ \mathcal{H}(\mathbf{y}, \mathbf{f}, \mathbf{n}, \boldsymbol{\sigma}) \triangleq [h_1, \dots, h_j, \dots, h_M]^T \\ h_j \triangleq \sigma_j \left[ \frac{C(n_j, y_j K / f_j)}{n_j f_j / K - y_j} + \frac{K}{f_j} \right] - \overline{D}_j. \end{cases}$$

Thus, MINLP **SP2** is

$$\begin{aligned} \min_{\mathbf{y} \in Y, \mathbf{f} \in F, \mathbf{n} \in N, \boldsymbol{\sigma} \in \Sigma} \quad & \mathcal{F}(\mathbf{f}, \mathbf{n}, \boldsymbol{\sigma}) \\ \text{s.t.} \quad & \begin{cases} \mathcal{G}(\mathbf{y}) = 0 \\ \mathcal{H}(\mathbf{y}, \mathbf{f}, \mathbf{n}, \boldsymbol{\sigma}) \leq \mathbf{0}. \end{cases} \end{aligned}$$

**Definition 3:** master problem  $\text{MP}^k$ :

$$\begin{aligned} \min_{\mathbf{n} \in N, \boldsymbol{\sigma} \in \Sigma, LB} \quad & LB \\ \text{s.t.} \quad & \begin{cases} LB \geq \mathcal{F}(\mathbf{f}^i, \mathbf{n}, \boldsymbol{\sigma}) + \lambda^i \mathcal{G}(\mathbf{y}^i) \\ \quad + (\boldsymbol{\mu}^i)^T \mathcal{H}(\mathbf{y}^i, \mathbf{f}^i, \mathbf{n}, \boldsymbol{\sigma}) \quad \forall i \in \mathcal{I}^k \\ 0 \geq \lambda^j \mathcal{G}(\mathbf{y}^j) + (\boldsymbol{\mu}^j)^T \mathcal{H}(\mathbf{y}^j, \mathbf{f}^j, \mathbf{n}, \boldsymbol{\sigma}) \quad \forall j \in \mathcal{J}^k \end{cases} \end{aligned} \quad (9a)$$

**Definition 4:** subproblem  $\text{SP}(\mathbf{n}^k, \boldsymbol{\sigma}^k)$ :

$$\begin{aligned} \min_{\mathbf{y} \in Y, \mathbf{f} \in F} \quad & \mathcal{F}(\mathbf{f}, \mathbf{n}^k, \boldsymbol{\sigma}^k) \\ \text{s.t.} \quad & \begin{cases} \mathcal{G}(\mathbf{y}) = 0 \\ \mathcal{H}(\mathbf{y}, \mathbf{f}, \mathbf{n}^k, \boldsymbol{\sigma}^k) \leq \mathbf{0}. \end{cases} \end{aligned}$$

**Definition 5:** subproblem feasibility-check  $\text{SPF}(\mathbf{n}^k, \boldsymbol{\sigma}^k)$ :

$$\begin{aligned} \min_{\mathbf{y} \in Y, \mathbf{f} \in F, \mathbf{s}} \quad & \mathbf{1}^T \mathbf{s} \\ \text{s.t.} \quad & \begin{cases} \mathcal{G}(\mathbf{y}) = 0 \\ \mathbf{s} \geq \mathcal{H}(\mathbf{y}, \mathbf{f}, \mathbf{n}^k, \boldsymbol{\sigma}^k). \end{cases} \end{aligned}$$

---

**Algorithm 1:** GBD Algorithm for solving **SP2**

---

```

/* Initialization */
1 Set  $k \leftarrow 1, \mathcal{I}^1 \leftarrow \emptyset, \mathcal{J}^1 \leftarrow \emptyset, UB^0 \leftarrow +\infty$ ;
2 while do
3   Solve  $\text{MP}^k$  by, e.g., branch and bound;
4   if feasible solution then
5     Obtain solution  $(\mathbf{n}^k, \boldsymbol{\sigma}^k, LB^k)$ ;
6   else if unbounded solution then
7     Choose arbitrary  $\mathbf{n}^k \in N$  and  $\boldsymbol{\sigma}^k \in \Sigma$ ;
8     Set  $LB^k \leftarrow -\infty$ ;
9   endif
10  Solve  $\text{SP}(\mathbf{n}^k, \boldsymbol{\sigma}^k)$  by, e.g., dual decomposition;
11  if feasible solution then
12    Obtain solution  $(\mathbf{y}^k, \mathbf{f}^k)$  and Lagrangian
13    multiplier  $(\lambda^k, \boldsymbol{\mu}^k)$ ;
14    Set  $UB^k \leftarrow \min \{UB^{k-1}, \mathcal{F}(\mathbf{f}^k, \mathbf{n}^k, \boldsymbol{\sigma}^k)\}$ ;
15    if  $|UB^k - LB^k| \leq \epsilon$  then /* Converged */
16      return  $(\mathbf{y}^k, \mathbf{f}^k, \mathbf{n}^k, \boldsymbol{\sigma}^k)$ ;
17    else /* Add feasible constraint */
18      Set  $\mathcal{I}^{k+1} \leftarrow \mathcal{I}^k \cup \{k\}, \mathcal{J}^{k+1} \leftarrow \mathcal{J}^k$ ;
19    endif
20  else if infeasible solution then
21    Solve  $\text{SPF}(\mathbf{n}^k, \boldsymbol{\sigma}^k)$  by, e.g., dual decomposition;
22    Obtain solution  $(\mathbf{y}^k, \mathbf{f}^k)$  and Lagrangian
23    multiplier  $(\lambda^k, \boldsymbol{\mu}^k)$ ;
24    Set  $UB^k \leftarrow UB^{k-1}$ ;
25    /* Add infeasible constraint */
26    Set  $\mathcal{I}^{k+1} \leftarrow \mathcal{I}^k, \mathcal{J}^{k+1} \leftarrow \mathcal{J}^k \cup \{k\}$ ;
27  endif
28  Set  $k \leftarrow k + 1$ ;
29 endwhile

```

---

## APPENDIX B

### SOLVE **SP3** USING HUNGARIAN ALGORITHM

The Hungarian algorithm is a combinatorial optimization approach that solves the assignment problem in polynomial time. We define

$$C_{ij} \triangleq \min \{l_i - x_i, \lambda_{ij}^{\max}, y_j\} \quad \forall i \in \mathcal{N}, j \in \mathcal{M}.$$

Thus, **SP3** can be equivalently transformed into a standard form of the assignment problem:

$$\begin{aligned} \min_{z_{ij}} \quad & \sum_{i \in \mathcal{N}} \sum_{j \in \mathcal{M}} d_{ij} C_{ij} z_{ij} \\ \text{s.t.} \quad & \begin{cases} \sum_{j \in \mathcal{M}} z_{ij} = 1 \quad \forall i \in \mathcal{N} \\ \sum_{i \in \mathcal{N}} z_{ij} = 1 \quad \forall j \in \mathcal{M} \\ 0 \leq z_{ij} \leq 1 \quad \forall i \in \mathcal{N}, j \in \mathcal{M}, \end{cases} \end{aligned} \quad (10)$$

where  $z_{ij}$  represents the assignment of the fog device  $i$  to the cloud server  $j$ , taking value 1 if the assignment is done and 0 otherwise. This formulation allows also fractional values, but there is always an optimal solution where the variables take integer values. This is because the constraint matrix is totally unimodular.

To illustrate the Hungarian algorithm for solving the above problem, without loss of generality, we consider a simple case with  $|\mathcal{N}|=4$  and  $|\mathcal{M}|=3$ . Since the two sets  $\mathcal{N}$  and  $\mathcal{M}$  should be of equal size, we add an additional dummy cloud server  $CS_3$ . The above problem can be viewed graphically: three fog devices  $FD_1, FD_2, FD_3$ , and  $FD_4$  as well as three cloud servers  $CS_1, CS_2, CS_3$ , and  $CS_4$  (including the dummy one). The lines from  $FD_i$  to  $CS_j$  represent the values of cost  $d_{ij}C_{ij}$ , with all  $d_{i4}C_{i4}$  setting to 0. For generality, we define the cost matrix to be the  $n \times n$  matrix:

$$C \triangleq \begin{bmatrix} d_{11}C_{11} & \dots & d_{1n}C_{1n} \\ \vdots & \ddots & \vdots \\ d_{n1}C_{n1} & \dots & d_{nn}C_{nn} \end{bmatrix}.$$

An assignment is a set of  $n$  entry positions in the cost matrix, no two of which lie in the same row or column. The sum of the  $n$  entries of an assignment is its cost. An assignment with the smallest possible cost is called an optimal assignment.

**Theorem 1:** If a number is added to or subtracted from all of the entries of any one row or column of a cost matrix, then one optimal assignment for the resulting cost matrix is also an optimal assignment for the original cost matrix.

---

**Algorithm 2:** Hungarian Algorithm for solving (10)

---

```

1 Subtract the smallest entry in each row from all the
  entries of its row;
2 Subtract the smallest entry in each column from all the
  entries of its column;
3 while do
4   Draw lines through appropriate rows and columns so
    that all the zero entries of the cost matrix are covered
    and the minimum number of such lines is used;
    /* Test for optimality */
5   if the minimum number of covering lines is  $n$  then
6     return an optimal assignment of zeros;
7   else if the minimum number of covering lines is less
    than  $n$  then /* An optimal assignment of
    zeros is not yet possible */
8     Determine the smallest entry not covered by any
    line;
9     Subtract this entry from each uncovered row;
10    Add this entry to each covered column;
11  endif
12 endw
```

---

**Algorithm 2** applies **Theorem 1** to a given  $n \times n$  cost matrix to find an optimal assignment. To illustrate **Algorithm 2** for solving (10), without loss of generality, we consider a simple case as shown in transformation (11). Step ① is to subtract 0 from each row. Step ② is to subtract 35 from column 1, 75 from column 2, 55 from column 3, and 0 from column

4. Step ③ is to cover all zeros with the minimum number of horizontal or vertical lines. Since the minimum number of covering lines is less than 4, we find that 10 is the smallest entry not covered by any line, and then subtract 10 from each uncovered row. Step ④ is to add 10 to each covered column. Since the minimum number of covering lines is 4, an optimal assignment of zeros is obtained. Step ⑤ is to make the same assignment for the original cost matrix. Thus, the optimal assignment for this case is  $z_{12}^*=z_{23}^*=z_{34}^*=z_{41}^*=1$  with the smallest cost of 175.

---

**Algorithm 3:** Update parameters in SP3

---

```

1 for  $i \in \mathcal{N}, j \in \mathcal{M}$  do
2   if  $z_{ij}^* == 1$  then
3     if  $C_{ij} == l_i - x_i$  then /* Remove  $i$  */
4        $\lambda_{ij}^* \leftarrow l_i - x_i$ ;
5        $\mathcal{N} \leftarrow \mathcal{N} \setminus \{i\}$ ;
6        $y_j \leftarrow y_j - \lambda_{ij}^*$ ;
7     else if  $C_{ij} == \lambda_{ij}^{\max}$  then /* Remove  $i \sim j$  */
8        $\lambda_{ij}^* \leftarrow \lambda_{ij}^{\max}$ ;
9        $l_i - x_i \leftarrow l_i - x_i - \lambda_{ij}^*$ ;
10       $y_j \leftarrow y_j - \lambda_{ij}^*$ ;
11       $C_{ij} \leftarrow \infty$ ;
12    else if  $C_{ij} == y_j$  then /* Remove  $j$  */
13       $\lambda_{ij}^* \leftarrow y_j$ ;
14       $l_i - x_i \leftarrow l_i - x_i - \lambda_{ij}^*$ ;
15       $\mathcal{M} \leftarrow \mathcal{M} \setminus \{j\}$ ;
16    endif
17  endif
18 endfor
```

---

Based on the optimal assignment to problem (10), we update the corresponding parameters in **SP3** according to **Algorithm 3**. Then we get a new assignment problem (10). In the same way, by adding additional dummy cloud servers, we have two sets of nodes with equal size, together with the corresponding cost matrix. Again, we apply Hungarian **Algorithm 2** to solve (10), obtain the optimal assignment, and update parameters in **SP3** according to **Algorithm 3**. This process goes so on and so forth until all the unprocessed requests have been dispatched from fog devices to cloud servers.

## REFERENCES

- [1] R. Deng, R. Lu, C. Lai, and T. H. Luan, "Towards power consumption-delay tradeoff by workload allocation in cloud-fog computing," in *Proc. IEEE ICC*, 2015, pp. 3909–3914.
- [2] R. Lu, H. Zhu, X. Liu, J. K. Liu, and J. Shao, "Toward efficient and privacy-preserving computing in big data era," *IEEE Network*, vol. 28, no. 4, pp. 46–50, 2014.
- [3] N. Kumar, S. Misra, J. Rodrigues, and M. Obaidat, "Coalition games for spatio-temporal big data in Internet of vehicles environment: a comparative analysis," *IEEE Internet of Things Journal*, vol. 2, no. 4, pp. 310–320, 2015.
- [4] C. Lai, R. Lu, D. Zheng, H. Li, and X. Shen, "Toward secure large-scale machine-to-machine communications in 3GPP networks: challenges and solutions," *IEEE Communications Magazine*, vol. 53, no. 12, pp. 12–19, 2015.

$$\begin{aligned}
 C &= \begin{bmatrix} 90 & 75 & 75 & 0 \\ 35 & 85 & 55 & 0 \\ 125 & 95 & 90 & 0 \\ 45 & 110 & 95 & 0 \end{bmatrix} \xrightarrow{\textcircled{1}} \begin{bmatrix} 90 & 75 & 75 & 0 \\ 35 & 85 & 55 & 0 \\ 125 & 95 & 90 & 0 \\ 45 & 110 & 95 & 0 \end{bmatrix} \xrightarrow{\textcircled{2}} \begin{bmatrix} \cancel{90} & \cancel{75} & \cancel{75} & \cancel{0} \\ \cancel{35} & \cancel{85} & \cancel{55} & \cancel{0} \\ \cancel{125} & \cancel{95} & \cancel{90} & \cancel{0} \\ \cancel{45} & \cancel{110} & \cancel{95} & \cancel{0} \end{bmatrix} \\
 &\xrightarrow{\textcircled{3}} \begin{bmatrix} 55 & 0 & 20 & 0 \\ 0 & 10 & 0 & 0 \\ 80 & 10 & 25 & -10 \\ 0 & 25 & 30 & -10 \end{bmatrix} \xrightarrow{\textcircled{4}} \begin{bmatrix} \cancel{55} & \cancel{0} & \cancel{20} & \cancel{0} \\ \cancel{0} & \cancel{10} & \cancel{0} & \cancel{0} \\ \cancel{80} & \cancel{10} & \cancel{25} & \cancel{-10} \\ \cancel{0} & \cancel{25} & \cancel{30} & \cancel{-10} \end{bmatrix} \xrightarrow{\textcircled{5}} \begin{bmatrix} 90 & 75 & 75 & 0 \\ 35 & 85 & 55 & 0 \\ 125 & 95 & 90 & 0 \\ 45 & 110 & 95 & 0 \end{bmatrix} \quad (11)
 \end{aligned}$$

- [5] T. H. Luan, L. X. Cai, J. Chen, X. Shen, and F. Bai, "Engineering a distributed infrastructure for large-scale cost-effective content dissemination over urban vehicular networks," *IEEE Transactions on Vehicular Technology*, vol. 63, no. 3, pp. 1419–1435, 2014.
- [6] S. He, J. Chen, X. Li, X. S. Shen, and Y. Sun, "Mobility and intruder prior information improving the barrier coverage of sparse sensor networks," *IEEE Transactions on Mobile Computing*, vol. 13, no. 6, pp. 1268–1282, 2014.
- [7] N. Lu, N. Cheng, N. Zhang, X. Shen, and J. W. Mark, "Connected vehicles: Solutions and challenges," *IEEE Internet of Things Journal*, vol. 1, no. 4, pp. 289–299, 2014.
- [8] The Network. Cisco Delivers Vision of Fog Computing to Accelerate Value from Billions of Connected Devices. [Online]. Available: <http://newsroom.cisco.com/press-release-content?articleId=1334100>
- [9] S. He, J. Chen, F. Jiang, D. K. Yau, G. Xing, and Y. Sun, "Energy provisioning in wireless rechargeable sensor networks," *IEEE Transactions on Mobile Computing*, vol. 12, no. 10, pp. 1931–1942, 2013.
- [10] F. Bonomi, R. Milito, P. Natarajan, and J. Zhu, "Fog computing: A platform for Internet of Things and analytics," in *Big Data and Internet of Things: A Roadmap for Smart Environments*. Springer, 2014, pp. 169–186.
- [11] R. Deng, Z. Yang, M.-Y. Chow, and J. Chen, "A survey on demand response in smart grids: Mathematical models and approaches," *IEEE Transactions on Industrial Informatics*, vol. 11, no. 3, pp. 570–582, 2015.
- [12] J. Chen, Q. Yu, B. Chai, Y. Sun, Y. Fan, and X. Shen, "Dynamic channel assignment for wireless sensor networks: A regret matching based approach," *IEEE Transactions on Parallel and Distributed Systems*, vol. 26, no. 1, pp. 95–106, 2015.
- [13] L. Atzori, A. Iera, and G. Morabito, "The Internet of Things: A survey," *Computer networks*, vol. 54, no. 15, pp. 2787–2805, 2010.
- [14] J. A. Stankovic, "Research directions for the Internet of Things," *Internet of Things Journal*, *IEEE*, vol. 1, no. 1, pp. 3–9, 2014.
- [15] A. Zanella, N. Bui, A. Castellani, L. Vangelista, and M. Zorzi, "Internet of Things for smart cities," *IEEE Internet of Things Journal*, vol. 1, no. 1, pp. 22–32, 2014.
- [16] S. K. Datta, C. Bonnet, and J. Haeri, "Fog computing architecture to enable consumer centric Internet of Things services," in *Proc. IEEE International Symposium on Consumer Electronics (ISCE)*, 2015, pp. 1–2.
- [17] T. H. Luan, L. Gao, Z. Li, Y. Xiang, and L. Sun, "Fog computing: Focusing on mobile users at the edge," *arXiv preprint arXiv:1502.01815*, 2015.
- [18] R. Suryawansh and G. Mandlik, "Focusing on mobile users at edge and Internet of Things using fog computing," *International Journal of Scientific Engineering and Technology Research*, vol. 4, no. 17, pp. 3225–3231, 2015.
- [19] L. Rao, X. Liu, L. Xie, and W. Liu, "Coordinated energy cost management of distributed Internet data centers in smart grid," *IEEE Transactions on Smart Grid*, vol. 3, no. 1, pp. 50–58, 2012.
- [20] L. Yu, T. Jiang, Y. Cao, and Q. Qi, "Carbon-aware energy cost minimization for distributed Internet data centers in smart microgrids," *IEEE Internet of Things Journal*, vol. 1, no. 3, pp. 255–264, 2014.
- [21] L. Rao, X. Liu, M. D. Ilic, and J. Liu, "Distributed coordination of Internet data centers under multi-regional electricity markets," *Proceedings of the IEEE*, vol. 100, no. 1, pp. 269–282, 2012.
- [22] S. Yi, C. Li, and Q. Li, "A survey of fog computing: Concepts, applications and issues," in *Proc. ACM Workshop on Mobile Big Data (Mobidata)*, 2015, pp. 37–42.
- [23] S. Yi, Z. Qin, and Q. Li, "Security and privacy issues of fog computing: A survey," in *Wireless Algorithms, Systems, and Applications*. Springer, 2015, pp. 685–695.
- [24] X. Wang, X. Chen, C. Yuen, W. Wu, and W. Wang, "To migrate or to wait: Delay-cost tradeoff for cloud data centers," in *Proc. IEEE Globecom*, 2014, pp. 2314–2319.
- [25] X. Wang, C. Yuen, N. U. Hassan, W. Wang, and T. Chen, "Migration-aware virtual machine placement for cloud data centers," in *Proc. IEEE ICC Workshop*, 2015, pp. 1940–1945.
- [26] M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. Katz, A. Konwinski, G. Lee, D. Patterson, A. Rabkin, I. Stoica *et al.*, "A view of cloud computing," *Communications of the ACM*, vol. 53, no. 4, pp. 50–58, 2010.
- [27] I. Stojmenovic and S. Wen, "The fog computing paradigm: Scenarios and security issues," in *Proc. Federated Conference on Computer Science and Information Systems (FedCSIS)*, 2014, pp. 1–8.
- [28] I. Stojmenovic, "Fog computing: A cloud to the ground support for smart things and machine-to-machine networks," in *Proc. Australasian Telecommunication Networks and Applications Conference (ATNAC)*, 2014, pp. 117–122.
- [29] X. Wang, C. Yuen, X. Chen, N. U. Hassan, and Y. Ouyang, "Cost-aware demand scheduling for delay tolerant applications," *Journal of Network and Computer Applications*, vol. 53, pp. 173–182, 2015.
- [30] M. Maleki, K. Dantu, and M. Pedram, "Power-aware source routing protocol for mobile ad hoc networks," in *Proc. ACM International Symposium on Low Power Electronics and Design*, 2002, pp. 72–75.
- [31] F. Ahmad and T. Vijaykumar, "Joint optimization of idle and cooling power in data centers while maintaining response time," in *ACM Sigplan Notices*, vol. 45, no. 3, 2010, pp. 243–256.
- [32] N. Gautam, *Analysis of Queues: Methods and Applications*. CRC Press, 2012.
- [33] R. Deng, Y. Zhang, S. He, J. Chen, and X. Shen, "Maximizing network utility of rechargeable sensor networks with spatiotemporally-coupled constraints," *IEEE Journal on Selected Areas in Communications*, DOI: 10.1109/JSAC.2016.2520181, to appear.
- [34] J. Ren, Y. Zhang, R. Deng, N. Zhang, D. Zhang, and X. Shen, "Joint channel access and sampling rate control in energy harvesting cognitive radio sensor networks," *IEEE Transactions on Emerging Topics in Computing*, DOI: 10.1109/TETC.2016.2555806, to appear.
- [35] S. P. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [36] Z. Yang, K. Long, P. You, and M.-Y. Chow, "Joint scheduling of large-scale appliances and batteries via distributed mixed optimization," *IEEE Transactions on Power Systems*, vol. 30, no. 4, pp. 2031–2040, 2015.
- [37] R. Deng, G. Xiao, and R. Lu, "Defending against false data injection attacks on power system state estimation," *IEEE Transactions on Industrial Informatics*, DOI: 10.1109/TII.2015.2470218, to appear.
- [38] D. Li and X. Sun, *Nonlinear Integer Programming*. Springer, 2006.
- [39] H. Kuhn, "The hungarian method for the assignment problem," *Naval Research Logistics*, vol. 52, no. 1, pp. 7–21, 2005.



**Ruilong Deng** (S'11-M'14) received the B.Sc. and Ph.D. degrees both in Control Science and Engineering from Zhejiang University, China, in 2009 and 2014, respectively.

He was a Visiting Scholar at Simula Research Laboratory, Norway, in 2011, and the University of Waterloo, Canada, from 2012 to 2013. He was a Research Fellow at Nanyang Technological University, Singapore, from 2014 to 2015. Currently, he is an AITF Postdoctoral Fellow with the Department of Electrical and Computer Engineering, University of

Alberta, Canada. His research interests include smart grid, cognitive radio, and wireless sensor network.

Dr. Deng currently serves as an Editor for *IEEE/KICS Journal of Communications and Networks*, and a Guest Editor for *IEEE TRANSACTIONS ON EMERGING TOPICS IN COMPUTING* and *Journal of Computer Networks and Communications*. He also serves/served as a Technical Program Committee Member for IEEE Globecom, IEEE ICC, IEEE SmartGridComm, EAI SGSC, etc.

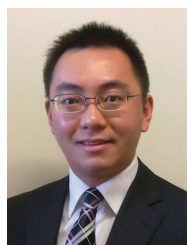


**Tom H. Luan** (M'13) received the B.Sc. degree from Xi'an Jiaotong University, China, in 2004, M.Phil. degree from Hong Kong University of Science and Technology in 2007, and Ph.D. degree from the University of Waterloo in 2012. Since December 2013, he has been the Lecturer in Mobile and Apps at the School of Information Technology, Deakin University, Melbourne Burwood, Australia. His research mainly focuses on vehicular networking, mobile content distribution, fog computing, and mobile cloud computing.



**Rongxing Lu** (S'09-M'11-SM'15) received the Ph.D. degree in computer science from Shanghai Jiao Tong University, Shanghai, China, in 2006, and the Ph.D. degree in electrical and computer engineering from the University of Waterloo, Waterloo, ON, Canada, in 2012. From May 2012 to April 2013, he was a Postdoctoral Fellow with the University of Waterloo. Since May 2013, he has been an Assistant Professor with the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore. His research interests include

computer network security, mobile and wireless communication security, and applied cryptography. Dr. Lu was the recipient of the Canada Governor General Gold Metal.



**Hao Liang** (S'09-M'14) is an Assistant Professor in the Department of Electrical and Computer Engineering at the University of Alberta, Canada, since 2014. He received his Ph.D. degree from the Department of Electrical and Computer Engineering, University of Waterloo, Canada, in 2013. From 2013 to 2014, he was a postdoctoral research fellow in the Broadband Communications Research (BBCR) Lab and Electricity Market Simulation and Optimization Lab (EMSOL) at the University of Waterloo. His current research interests are in the areas of smart grid, wireless communications, and wireless networking. He is a recipient of the Best Student Paper Award from IEEE 72nd Vehicular Technology Conference (VTC Fall-2010), Ottawa, ON, Canada.

Dr. Liang serves/served as a Guest Editor for *IEEE TRANSACTIONS ON EMERGING TOPICS IN COMPUTING* and *Journal of Computer Networks and Communications*. He has been a Technical Program Committee (TPC) Member for major international conferences in both information/communication system discipline and power/energy system discipline, including IEEE International Conference on Communications (ICC), IEEE Global Communications Conference (Globecom), IEEE VTC, IEEE Innovative Smart Grid Technologies Conference (ISGT), and IEEE International Conference on Smart Grid Communications (SmartGridComm). He was the System Administrator of *IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY* (2009-2013).



**Chengzhe Lai** (M'14) received his degree in B.S. in Information Security from Xi'an University of Posts and Telecommunications in 2008 and a Ph.D. degree from Xidian University in 2014. At present, he is with the School of Telecommunication and Information Engineering, Xi'an University of Posts and Telecommunications and with the National Engineering Laboratory for Wireless Security, Xi'an, China. His research interests include wireless network security, privacy preservation, and M2M communications security.