

# PAM & PAL: Policy-Aware Virtual Machine Migration and Placement in Dynamic Cloud Data Centers

Hugo Flores, Vincent Tran, and Bin Tang

Department of Computer Science

California State University Dominguez Hills, Carson, CA 90747, USA

Email: {hflores27,vtran42}@toromail.csudh.edu, btang@csudh.edu

**Abstract**—We focus on policy-aware data centers (PADCs), wherein virtual machine (VM) traffic traverses a sequence of middleboxes (MBs) for security and performance purposes, and propose two new VM placement and migration problems. We first study PAL: policy-aware virtual machine placement. Given a PADC with a data center policy that communicating VM pairs must satisfy, the goal of PAL is to place the VMs into the PADC to minimize their total communication cost. Due to dynamic traffic loads in PADCs, however, above VM placement may no longer be optimal after some time. We thus study PAM: policy-aware virtual machine migration. Given an existing VM placement in the PADC and dynamic traffic rates among communicating VMs, PAM migrates VMs in order to minimize the total cost of migration and communication of the VM pairs. We design optimal, approximation, and heuristic policy-aware VM placement and migration algorithms. Our experiments show that i) VM migration is an effective technique, reducing total communication cost of VM pairs by 25%, ii) our PAL algorithms outperform state-of-the-art VM placement algorithm that is oblivious to data center policies by 40-50%, and iii) our PAM algorithms outperform the only existing policy-aware VM migration scheme by 30%.

**Index Terms**—Policy-Aware Data Centers, Virtual Machine Placement, Virtual Machine Migration, Algorithms

## I. INTRODUCTION

Recently, middleboxes (MBs) [10], such as firewalls, load balancers, and network address translators, are introduced into cloud data centers to improve security and performances of virtual machine (VM) applications [40]. In particular, *data center policies* (or *service function chaining*) are established that require VM traffic to traverse a chain of MBs [26], [30], [45]. Fig. 1(a) shows such an example. Traffic generated at VM  $vm_1$  goes through a firewall, a load balancer, and a cache proxy before arriving at VM  $vm'_1$ . In doing so, this policy filters out malicious traffic, diverts trusted VM traffic to avoid network congestion, and finally caches the content to share with other cloud users. We refer to cloud data centers that implement such policies as *policy-aware data centers* (PADCs). Fig. 1(b) shows a small PADC that implements the same data center policy in Fig. 1(a). A firewall, a load balancer, and a cache proxy are installed on switches  $s_2$ ,  $s_3$ , and  $s_4$ , respectively. There are four physical machines (PMs),  $pm_1$ ,  $pm_2$ ,  $pm_3$ , and  $pm_4$ , each can store one VM due to its limited *cloud resources* (i.e., CPU, memory, and disk I/O).

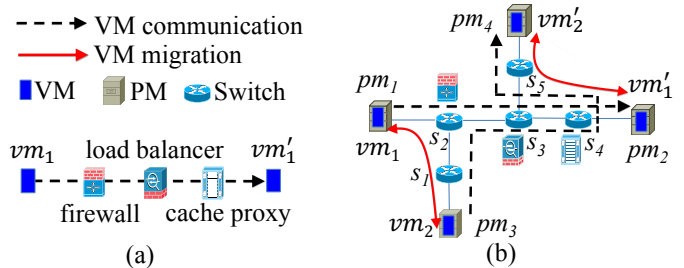


Fig. 1. (a) A data center policy. (b) A PADC example.

We identify, formulate, and solve two new VM placement and migration problems in PADCs. Measurements from Facebook and other production data centers show that traffic loads (i.e., transmission rates) of communicating VM applications are highly diverse and dynamic [9], [36]. One such example is cloud chatbots (e.g., Slack [4] and Amazon Lex [1]), wherein VMs could send to each other low-bandwidth texts at one instant then switch to high-bandwidth live-streaming videos at another and vice versa, while the duration of each communication varies considerably. Thus it is important to design resource-efficient VM placement and migration for dynamic cloud data centers. This is especially crucial for PADCs – as the VM communication must go through a sequence of MBs, it generates more network traffic and consumes more *network resources* such as bandwidths compared to traditional cloud data centers. In this paper we focus on pairwise VM communication as the east-west traffic within data centers accounts for 75.4 percent of data center traffic [2], and most east-west traffic in cloud data centers is pairwise [33].

Initially, when new cloud applications are launched as VMs in PADCs, it needs to decide how to place them for resource-efficient communication. We refer to this problem as PAL: *policy-aware VM placement in PADCs*. Given a PADC with a data center policy that communicating VM pairs must satisfy, PAL studies how to place the VMs with diverse traffic rates into the PMs to minimize their total communication cost while satisfying resource constraints of PMs. Fig. 1(b) shows that two VM pairs ( $vm_1, vm'_1$ ) and ( $vm_2, vm'_2$ ), with the traffic load of the former much larger than that of the latter, are to be placed. To reduce network traffic and communication delay, it is preferred that  $vm_1$  and  $vm'_1$  communicates along

a route that is “shorter” than that of  $vm_2$  and  $vm'_2$ . One solution is to place  $vm_1$  and  $vm'_1$  on  $pm_1$  and  $pm_2$ , and  $vm_2$  and  $vm'_2$  on  $pm_3$  and  $pm_4$ , with their respective policy-aware communication routes shown in black dashed lines.

However, such initial VM placement may become suboptimal due to dynamic traffic in PADCs. In Fig. 1(b), if the traffic load of  $(vm_2, vm'_2)$  emerges as much larger than that of  $(vm_1, vm'_1)$ , the aforesaid VM placement becomes inefficient – as  $vm_2$  communicates with  $vm'_2$  via a route much longer than that of  $(vm_1, vm'_1)$ , it generates more network traffic and consumes much of the network bandwidths. To tackle such problem, we observe that migrating VMs from one PM to another might be an effective technique. In Fig. 1(b), it might be a good idea to “swap” these two VM pairs using VM migration routes shown in red solid lines. Given an existing placement of communicating VM pairs with dynamic traffic rates, a data center policy that they must satisfy, we study how to migrate VMs to minimize the total cost of migration and communication of all the VM pairs. We refer to this problem as PAM: *policy-aware VM migration in PADCs*.

Consider that VM migration incurs traffic overhead, and that a large scale PADC typically has tens of thousands of PMs, as well as hundreds of thousands of communicating VMs with wide range of changing traffic rates, how to effectively place and migrate VMs in PADCs becomes challenging problem. In this paper we address this challenge by formulating PAL and PAM and design optimal, approximation, and heuristic *policy-aware* algorithms to solve them. The PAM & PAL duo potentially achieves ideal resource utilization for a PADC’s lifetime - after the PAL algorithms create the initial VM placement to optimize a PADC’s cloud resource utilization, the PAM algorithms can then be executed periodically to optimize a PADC’s network resource utilization while adapting dynamic VM traffic. To the extent of our knowledge, PAL and PAM are new problems that have not been studied before.

Using traffic patterns and flow characteristics found in production data centers, we show that VM migration reduces the total communication costs of VM pairs by 25%, demonstrating that it is an effective technique to alleviate dynamic VM traffic in PADCs. We also show our *policy-aware* algorithms outperform the *only* existing policy-aware VM migration algorithm [15] by 30%, and the state-of-the-art VM placement algorithm that is oblivious to data center policies [33] by 40-50%.

## II. RELATED WORK

Service function chaining (SFC) has been an active research in recent years. It mainly focused on virtual network functions (VNFs) (i.e., virtual MBs) with their implementation and realization [35], [19], [20], [45], [25], VNF placements [27], [34], VNF migrations [17], [31], and other issues such as availability [18], flow control [38], and finding shortest SFC [37]. However, given that virtual MBs cannot yet match the performance of hardware MBs thus many network functions still rely on dedicated hardware, and much existence of hardware MBs in enterprise networks [24], [39], we consider hardware

MBs in this paper. As such, once they are installed inside data centers, they cannot be easily moved around.

There is vast amount of literature of VM migration in cloud data centers [41], [44], [8], [42], [43], [16]. In particular, Shrivastava et al. [41] proposed an application-aware VM migration that minimizes data center network traffic while considering the combined effects of application dependencies and network topology. Zhang et al. [44] analyzed how much bandwidth is required to guarantee the total migration time and downtime of a live VM migration. Wang et al. [42] studied how to migrate multiple VMs at the same time with available bandwidth, and designed a fully polynomial time approximation algorithm. Cui et al. [16] assumed that data center topologies are adaptive with reconfigurable wireless links or optical circuit switches, and proposed VM migration algorithm with constant approximation ratio.

VM migration studied in this paper, however, differs from aforesaid work in both goals and models. While existing VM migration work achieved various objectives such as server consolidation and energy efficiency, load balancing and fault tolerance, our work focuses on the dynamic communication traffic rates existing among VMs. Besides, none of the above work considered data center policies, thus falling short of achieving performance and security guarantees brought about by various of MBs deployed inside PADCs.

Meng et al. [33] designed one of the first *policy-oblivious* VM placement algorithms. It is traffic-aware in that it assigns VMs with large communications to the same PMs or PMs in close proximity. As PAL is the first to study *policy-aware* VM placement thus no closely related work to compare with, we compare our work with this traffic-aware VM placement. Alicherry and Lakshman [7] designed optimal and approximation algorithms that place VMs to minimize data access latencies. Li et al. [29] studied the VM placement to reduce the data center network cost as well as the cost utilizing PMs. Again, they are policy-oblivious thus do not achieve performance and security guarantees brought by PADCs.

PACE [28] was one of the first to study *policy-aware* VM placement. However, it only considers one type of MB thus does not study data center policy addressed in this paper. The only closely related work to ours is by Cui et al. [15], [14]. They proposed PLAN, the first policy-aware VM migration scheme for all-pair VM communications, and provided heuristic algorithms for ordered policies. In contrast, we focus on pairwise VM communication. We design optimal, approximate, and heuristic algorithms for both ordered and unordered policies and show they all outperform PLAN.

One salient feature of our research is that it integrates VM migration and VM placement, two most significant VM management mechanisms, into one framework. Existing VM migration and placement research seems achieving disparate goals – While VM placement optimizes *communication costs* among VMs such as energy cost, data access delay, and bandwidth usage [7], [33], [29], [12], VM migration optimizes *migration costs* of VMs including migration time, downtime, and service degradation during migration [42], [44], [16]. By

modeling *topology-aware* VM migration and communication costs, we are able to jointly optimize VM placement and migration for the overall resource utilization in PADCs. We believe our work is the first that takes a holistic approach to solve VM placement and migration in PADCs.

### III. SYSTEM MODELS

**Network Model.** We use fat tree [6] to illustrate the problems and their algorithmic solutions. However, as the problems are applicable to any data center topology, we model a PADC as an undirected general graph  $G(V, E)$ .  $V = V_p \cup V_s$  is a set of PMs  $V_p = \{pm_1, pm_2, \dots, pm_{|V_p|}\}$  and a set of switches  $V_s$ .  $E$  is a set of edges, each connecting either one switch to another switch or a switch to a PM. Fig. 2 shows a  $k = 4$  PADC where  $k$  is the number of ports each switch has.

There are a set of  $n$  MBs  $\mathcal{M} = \{mb_1, mb_2, \dots, mb_n\}$  inside the PADC, with  $mb_j$  installed at switch  $sw_j \in V_s$ . We adopt the *bump-off-the-wire* design [26], which uses a policy-aware switching layer to explicitly redirect traffic to off-path MBs. Fig. 2 shows three MBs  $mb_1$ ,  $mb_2$  and  $mb_3$  installed using this design. As a switch and its attached MB are connected by high-speed optical fibers, the delay between them is negligible compared to that among switches and PMs [22].

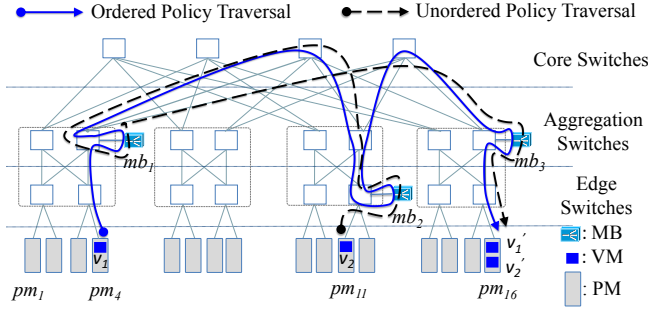


Fig. 2. A PADC with 16 PMs:  $pm_1, pm_2, \dots$ , and  $pm_{16}$ , 3 MBs:  $mb_1, mb_2$ , and  $mb_3$ , and two VM pairs:  $(v_1, v'_1)$  and  $(v_2, v'_2)$ .  $\bullet$  and  $\blacktriangleright$  indicate source and destination VM respectively.

There are  $l$  pairs of communicating VMs  $\mathcal{P} = \{(v_1, v'_1), (v_2, v'_2), \dots, (v_l, v'_l)\}$  that are already placed into the PMs. For any VM pair  $(v_i, v'_i)$ ,  $1 \leq i \leq l$ ,  $v_i$  and  $v'_i$  are referred to as its *source* and *destination VM* respectively. Denote the traffic rate or transmission rate of  $(v_i, v'_i)$  as  $\lambda_i$ , and the *traffic rate vector* as  $\vec{\lambda} = \langle \lambda_1, \lambda_2, \dots, \lambda_l \rangle$ . In a dynamic PADC,  $\vec{\lambda}$  is not a constant vector as the VM traffic rates change over time. In Fig. 2, there are two VM pairs:  $(v_1, v'_1)$  and  $(v_2, v'_2)$ , with  $\vec{\lambda} = \langle 100, 1 \rangle$ .

Let  $\mathcal{V} = \{v_1, v'_1, v_2, v'_2, \dots, v_l, v'_l\}$ . We assume that it needs one unit of cloud resource to create and execute a VM in  $\mathcal{V}$  and leave the more general case with varying resource demands as future work. Here the *resource* is an aggregated characterization of a PM's hardware resources such as CPUs, memories, and disk I/O. Denote the *resource capacity* of PM  $pm_i$  as  $rc_i$ , meaning that  $pm_i$  has  $rc_i$  resource slots. As there are  $2 \cdot l$  VMs and each needs one resource slot, it must be that  $\sum_{i=1}^{|V_p|} rc_i \geq 2 \cdot l$ . Table I shows all the notations.

TABLE I  
NOTATION SUMMARY

Notation	Description
$V_p$	$V_p = \{pm_1, pm_2, \dots, pm_{ V_p }\}$ is the set of $ V_p $ PMs
$V_s$	Set of switches in a PADC
$\mathcal{M}$	$\mathcal{M} = \{mb_1, mb_2, \dots, mb_n\}$ is the set of $n$ MBs
$\mathcal{P}$	$\mathcal{P} = \{(v_i, v'_i), \dots, (v_l, v'_l)\}$ is the set of $l$ VM pairs
$\mathcal{V}$	$\mathcal{V} = \{v_1, \dots, v_l, v'_1, \dots, v'_l\}$
$\lambda_i$	Traffic rate between $v_i$ and $v'_i$ , $1 \leq i \leq l$
$\vec{\lambda}$	$\vec{\lambda} = \langle \lambda_1, \lambda_2, \dots, \lambda_l \rangle$
$rc_i$	Resource capacity of PM $pm_i$ , $1 \leq i \leq  V_p $
$sw_j$	Switch where $mb_j$ is installed, $1 \leq j \leq n$
$c(i, j)$	Communication cost between PMs (or switches) $i$ and $j$
$p(v)$	PM where the VM $v$ is placed with VM placement $p$
$\pi^i$	Order at which $(v_i, v'_i)$ visits MBs in unordered policy
$\vec{\pi}$	$\vec{\pi} = \langle \pi^1, \pi^2, \dots, \pi^l \rangle$
$C_c(p)$	Total communication cost with $p$ in ordered policy
$C_c(p, \vec{\pi})$	Total communication cost with $p$ in unordered policy
$\mu$	Migration coefficient
$m(v)$	PM where the VM $v$ migrates to under VM migration $m$
$C_m(m)$	Total migration cost with migration $m$
$C_c(m)$	Total communication cost after migration $m$
$C_t(m)$	Total migra. and comm. cost with $m$ in ordered policy
$C_t(m, \vec{\pi})$	Total migra. and comm. cost with $m$ in unordered policy

**Cost Model.** We model the VM communication and migration cost as either the delay or energy cost of the network traffic inside PADCs. We adopt a *topology-aware* model [33] and define the *communication cost* of any VM pair as the number of network links its traffic traverses multiplied by its traffic rate (however, our problems and solutions still hold for that different edges have different costs). The *migration cost* of migrating any VM  $v$  from PM  $i$  to PM  $j$  is  $\mu \cdot c(i, j)$ . Here,  $c(i, j)$  is the minimum number of hops between any PM (or switch)  $i$  and  $j$ , and  $\mu$  is a *migration coefficient* that depends on VM sizes and available network bandwidths.

**Justifications.** Our VM migration model is different from the one adopted by most existing literature. Mann et al. [32] focused on *pre-copy* [11], one of the original live VM migration techniques, and modeled the cost of migrating a VM  $v$  as  $M_s \cdot \frac{1 - (P_r/B_a)^{n+1}}{1 - (P_r/B_a)}$ . Here  $M_s$  is the image size of  $v$ ,  $P_r$  is its page dirty rate,  $B_a$  is the available bandwidth, and  $n$  is number of pre-copy phases. They suggested the migration cost be a constant quantity measured by the hypervisor. In contrast, by acknowledging network delay or energy consumption incurred by VM migration traffic, our topology-aware model is more conducive to designing VM migration algorithms for a large-scale and dynamical traffic environment targeted by this paper.

**Data Center Policies.** Depending on the application requirements, some policies require that the VM traffic to go through the MBs in a strict order. We refer to such policies as *ordered policies* and denote them as  $(mb_1, mb_2, \dots, mb_n)$ . On the other hand, as MB functions are mostly independent from each other, some data center policies are considered satisfied as long as all its MBs are visited by the VM traffic. We refer to such policies as *unordered policies* and denote them as  $\{mb_1, mb_2, \dots, mb_n\}$ . In Fig. 2,  $(v_1, v'_1)$  traverses MBs under ordered policy  $(mb_1, mb_2, mb_3)$ , resulting in communication cost of  $100 \times 10 = 1000$  (solid blue line),  $(v_2, v'_2)$  traverses MBs under unordered policy  $\{mb_1, mb_2, mb_3\}$ , resulting in communication cost of  $1 \times 8 = 8$  (dashed black line).

We refer to the first (and last) visited MB in a policy as *ingress (and egress) MB*, and the switch where the ingress (and egress) MB is installed as *ingress (and egress) switch*. For ordered policy, the ingress switch is  $sw_1$  and egress switch is  $sw_n$ . For unordered policy, it needs to find out  $sw_1$  and  $sw_n$  as well as the MB sequence along which VM pair communicates. As one data center policy as the one shown in Fig. 1(a) is generally sufficient to serve both security and performance purposes, we assume there is one data center policy (ordered or unordered) in a PDDC at a time. We adopt FlowTags [19], a SDN architecture that provides consistent policy enforcement during VM migration by adding tags in packet headers.

**EXAMPLE 1:** Fig. 3 shows a  $k = 2$  linear fat tree PADC with two PMs:  $pm_1$  and  $pm_2$ . Each PM has two resource slots; the four of them are  $\{rs_1, rs_2, rs_3, rs_4\}$ . Two MBs,  $mb_1$  and  $mb_2$ , are installed on edge switch  $sw_1$  and aggregation switch  $sw_2$ , respectively. There are two VM pairs  $(v_1, v'_1)$  and  $(v_2, v'_2)$ ,  $v_1$  and  $v_2$  are placed on  $pm_1$  while  $v'_1$  and  $v'_2$  on  $pm_2$ .  $\bar{\lambda} = \langle 100, 1 \rangle$  and  $\mu = 1$ . Before migration, the total communication cost in Fig. 3(a) is 606 under both  $(mb_1, mb_2)$  and  $\{mb_1, mb_2\}$ . By migrating  $v'_1$  to  $pm_1$  and  $v_2$  to  $pm_2$  with migration cost of 12 (solid red line in Fig. 3(a)), the total communication cost (dotted and dashed black lines in Fig. 3(b)) becomes 410, a 30% of total cost reduction. We show this migration is indeed optimal in Section V-A.  $\square$

#### IV. PAL: POLICY-AWARE VM PLACEMENT IN PADCs

##### A. Ordered Policy.

1) *Problem Formulation:* Under ordered policy, for any VM pair communication, the ingress switch is always  $sw_1$  and the egress switch is always  $sw_n$ . Given a VM placement function  $p$ , denote the *total communication cost* of all the  $l$  VM pairs under  $p$  as  $C_c(p)$ . We have  $C_c(p) =$

$$= \sum_{i=1}^l \lambda_i \cdot \sum_{j=1}^{n-1} c(sw_j, sw_{j+1}) + \sum_{i=1}^l \lambda_i \cdot \left( c(p(v_i), sw_1) + c(sw_n, p(v'_i)) \right). \quad (1)$$

The objective of PAL is to find a VM placement  $p$  to minimize  $C_c(p)$  while satisfying resource constraint of PMs:  $|\{v \in \mathcal{V} | p(v) = i\}| \leq rc_i, \forall i \in V_p$ . As the first term on the r.h.s. of Eq. 1 is fixed under ordered policy, we only need to minimize the second term. Below we design an optimal and efficient algorithm to solve PAL.

2) *VM Placement Algorithm for Ordered Policy:* To save communication costs for VM pairs, the key is to find a set of resource slots close to the ingress (and egress) switch to place source (and destination) VMs. We give below definitions.

**Definition 1: (Ingress/Egress Costs, Ingress/Egress Resource Sets, Optimal Ingress/Egress Sets)** A resource slot  $rs$ 's *ingress* (and *egress*) cost, denoted as  $c_{in}(rs)$  (and  $c_e(rs)$ ), is the cost between its belonged PM and the ingress switch  $sw_1$  (and egress switch  $sw_n$ ). Let  $pm(rs)$  be the PM  $rs$  belongs to,  $c_{in}(rs) = c(pm(rs), sw_1)$ ,  $c_e(rs) = c(pm(rs), sw_n)$ .

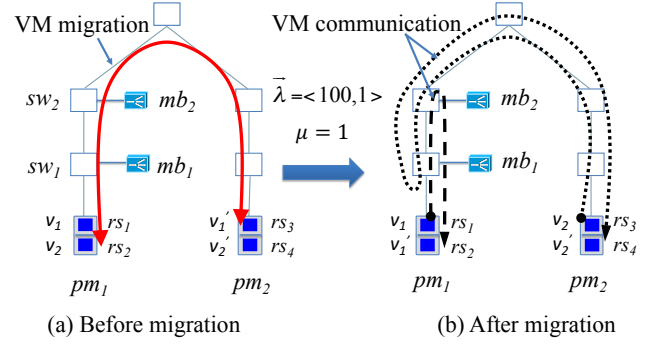


Fig. 3. VM migration achieved 30% of total cost reduction in a linear PADC.

An *ingress (and egress) resource set* (IRS and ERS) is a set of  $l$  resource slots that store the  $l$  source (and destination) VMs. The *cost* of an IRS (and ERS) is the sum of the ingress (and egress) costs of its resource slots. A pair of IRS and ERS is *optimal*, denoted as  $(I^{opt}, E^{opt})$ , if the sum of their costs is the minimum among all pairs of IRS and ERS.  $\square$

$I^{opt}$  and  $E^{opt}$  are structures that uniquely arise in PAL. Algo. 1 below finds such a pair (lines 1-22) and then places the  $l$  VM pairs (in non-ascending order of their traffic rates) into it (lines 23-30). Its time complexity is  $O(|V_p|^2 \cdot \bar{m}^2)$ , where  $\bar{m}$  is the average resource capacity of a PM.

##### Algorithm 1: PAL Algorithm for Ordered Policy.

**Input:** A PADC with ordered policy  $(mb_1, mb_2, \dots, mb_n)$ ,

VM pairs  $\mathcal{P}$ ,  $V_p = \{pm_i\}$ , resource capacity  $rc_i$ .

**Output:** A placement  $p$  and the total comm. cost  $C_c(p)$ .

**Notations:**  $\mathcal{I}$  and  $\mathcal{E}$ : arrays of resource slots, each of size  $2l$ .  $i, j$ : indices for  $\mathcal{I}$  and  $\mathcal{E}$  respectively.

$k$ : index for  $I^{opt}$  and  $E^{opt}$ .

$sel(rs_i)$ : initially *false*, *true* if  $rs_i$  is put into  $I^{opt}$  or  $E^{opt}$ .

1.  $i = j = k = 1$ ,  $C_c(p) = 0$ ,  $p = \phi$ ;
2. Sort resource slots in non-descending order of their ingress (and egress) costs, store the top  $2l$  in arrays  $\mathcal{I}$  (and  $\mathcal{E}$ );
3. **while** ( $k \leq l$ ) // find optimal resource slots for  $(v_k, v'_k)$
4. **if** ( $sel[\mathcal{I}[i]] == true$ )  $i++$ ;
5. **if** ( $sel[\mathcal{E}[j]] == true$ )  $j++$ ;
6. **if** ( $\mathcal{I}[i] \neq \mathcal{E}[j]$ ) // both optimal resource slots are found
7.  $I^{opt}[k] = \mathcal{I}[i]$ ,  $E^{opt}[k] = \mathcal{E}[j]$ ;
8.  $sel[\mathcal{I}[i]] = sel[\mathcal{E}[j]] = true$ ;
9.  $i++, j++$ ;
10. **else** // one found, now find the other
11. **if** ( $c_{in}(\mathcal{I}[i]) + c_e(\mathcal{E}[j+1]) \leq c_{in}(\mathcal{I}[i+1]) + c_e(\mathcal{E}[j])$ )
12.  $I^{opt}[k] = \mathcal{I}[i]$ ,  $E^{opt}[k] = \mathcal{E}[j+1]$ ;
13.  $sel[\mathcal{I}[i]] = sel[\mathcal{E}[j+1]] = true$ ;
14.  $i++, j += 2$ ;
15. **else**
16.  $I^{opt}[k] = \mathcal{I}[i+1]$ ,  $E^{opt}[k] = \mathcal{E}[j]$ ;
17.  $sel[\mathcal{I}[i+1]] = sel[\mathcal{E}[j]] = true$ ;
18.  $i += 2, j++$ ;
19. **end if**;
20. **end if**;
21.  $k++$ ;
22. **end while**;
23. **WLOG**,  $\lambda_1 \geq \lambda_2 \dots \geq \lambda_l$ ;

24. **for** ( $1 \leq i \leq l$ ) // place VM pairs and calculate cost  
25. Place  $v_i$  at  $I^{opt}[i]$ ,  $v'_i$  at  $E^{opt}[i]$ ;  
26.  $p = p \cup \{(I^{opt}[i], E^{opt}[i])\}$ ;  
27.  $C_c(p) += \lambda_i * (c_{in}(I^{opt}[i]) + c_e(E^{opt}[i]))$ ;  
28. **end for**;  
29.  $C_c(p) += \sum_{i=1}^l \lambda_i \sum_{j=1}^{n-1} c(sw_j, sw_{j+1})$ ;  
30. **RETURN**  $p$  and  $C_c(p)$ .

**EXAMPLE 2:** Fig. 4(a) shows how Algo. 1 could place the two VM pairs  $(v_1, v'_1)$  and  $(v_2, v'_2)$  into the same PADC in Fig. 3. It gives  $\mathcal{I} = \mathcal{E} = \{rs_1, rs_2, rs_3, rs_4\}$ , from which it computes  $I^{opt} = \{rs_1, rs_3\}$  and  $E^{opt} = \{rs_2, rs_4\}$ . It thus places  $v_1$  and  $v'_1$  in  $pm_1$  and  $v_2$  and  $v'_2$  in  $pm_2$  with total communication cost of  $100 \cdot 4 + 1 \cdot 10 = 410$ .  $\square$

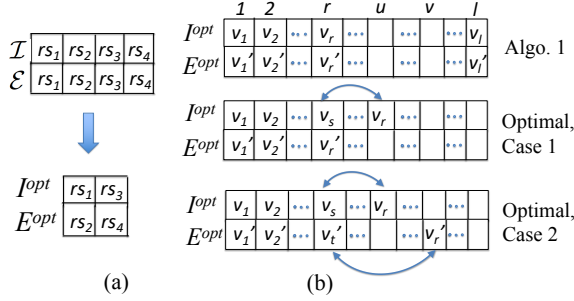


Fig. 4. (a) A working example and (b) optimality proof for Algo. 1.

**Theorem 1:** Algo. 1 finds the VM placement that minimizes total communication cost for the  $l$  VM pairs.

**Proof Sketch.** First, we prove by induction that  $(I^{opt}, E^{opt})$  computed in Algo. 1 (lines 1-22) is a pair of optimal IRS and ERS. Second, we prove by contradiction that the VM placement on  $I^{opt}$  and  $E^{opt}$  in Algo. 1 (lines 23-30) yields minimum total communication cost. Assume that Algo. 1 is not optimal and that  $r$ ,  $1 \leq r \leq l$ , is the smallest index at which  $I^{opt}[r]$  or  $E^{opt}[r]$  store different pair of VMs for Algo. 1 and Optimal. Two cases are shown in Fig. 4(b). Case 1: one of the two resource slots,  $I^{opt}[r]$  or  $E^{opt}[r]$ , stores different VMs. For example, both algorithms stores  $v'_r$  at  $E^{opt}[i]$  while Algo. 1 stores  $v_r$  and Optimal stores  $v_s$  at  $I^{opt}[i]$ . Case 2: both slots store different VMs. For example, Optimal stores  $v_s$  at  $I^{opt}[r]$  and  $v'_t$  at  $E^{opt}[r]$  (with  $s, t > r$ ), and stores  $v_r$  at  $I^{opt}[u]$  and  $v'_r$  at  $E^{opt}[v]$  (with  $u, v > r$ ). In both cases, we can swap VMs in Optimal following blue curved arrow lines to further reduce its cost, due to  $\lambda_1 \geq \lambda_2 \dots \geq \lambda_l$ .  $\blacksquare$

### B. Unordered Policy.

1) *Problem Formulation.*: In unordered policy, besides a VM placement function  $p$ , PAL needs to find the order for each VM pair to visit the MBs. For  $(v_i, v'_i)$  we define such order as an MB traversal function  $\pi^i : [1, 2, \dots, n] \rightarrow [1, 2, \dots, n]$ , a permutation function indicating the  $j^{th}$  MB that  $(v_i, v'_i)$  visits is  $mb_{\pi^i(j)}$ . Given  $p$  and  $\pi^i$ , denote  $(v_i, v'_i)$ 's communication cost as  $c_{i, \pi^i}^{p, \pi^i}$ . Then  $c_{i, \pi^i}^{p, \pi^i} = \lambda_i \cdot c(p(v_i), sw(\pi^i(1))) + \lambda_i \cdot \sum_{j=1}^{n-1} c(sw(\pi^i(j)), sw(\pi^i(j+1))) + \lambda_i \cdot c(sw(\pi^i(n)), p(v'_i))$ . Let  $\vec{\pi} = \langle \pi^1, \pi^2, \dots, \pi^l \rangle$ . The objective of PAL under unordered policy is to minimize total

communication cost  $C_c(p, \vec{\pi}) = \sum_{i=1}^l c_{i, \pi^i}^{p, \pi^i}$  while satisfying  $|\{v \in \mathcal{V} | p(v) = i\}| \leq rc_i, \forall i \in V_p$ . Below we show that PAL is NP-hard even for one pair of VMs. We then propose an approximation algorithm for this special case that yields total cost at most twice of the optimal.

**Theorem 2:** Under unordered policy, PAL is NP-hard even for one pair of VMs  $(v_1, v'_1)$  (i.e.,  $l = 1$ ).

**Proof:** We reduce *s-t traveling salesman path problem* (TSPP) [23], which is NP-hard, to this problem. Given a complete undirected graph  $K = (V_K, E_K)$  with edge cost  $d : E_K \rightarrow \mathbb{R}^+$  and a pair of pre-specified vertices  $s, t \in V_K$ , TSPP finds a shortest *Hamiltonian path* that starts at  $s$ , visits each vertex exactly once, and ends at  $t$ .  $d$  satisfies *triangle inequality*, i.e.,  $d(u, v) \leq d(u, w) + d(w, v)$  for all  $u, v, w \in V_K$ . When  $s = t$ , TSPP becomes well-known *traveling salesman problem* (TSP) [13], which finds a shortest *Hamiltonian cycle*.

Given VM pair  $(v_1, v'_1)$  and an instance of PADC graph  $G(V = V_p \cup V_s, E)$ , we construct  $|V_p| \cdot (|V_p| + 1)/2$  instances of complete graphs  $K^{i,j} = (V_K^{i,j}, E_K^{i,j})$ ,  $1 \leq i \leq |V_p|$ ,  $i \leq j \leq |V_p|$ . Here,  $V_K^{i,j} = \{pm_i, pm_j, sw_1, sw_2, \dots, sw_n\}$  and for  $(u, v) \in E_K^{i,j}$ , its cost  $d(u, v)$  is the cost of the shortest path connecting  $u$  and  $v$  in  $G$ . Now, if  $K^{a,b}$  has a shortest Hamiltonian path whose cost is the minimum among the shortest Hamiltonian paths in all the instances, then placing  $v_1$  to  $pm_a$  and  $v'_1$  to  $pm_b$  must give the minimum communication cost for  $(v_1, v'_1)$  in  $G$ , and vice versa.  $\blacksquare$

### 2) VM Placement Algorithm for Unordered Policy:

**Definition 2: (Optimal Policy Route (OPR).)** In a PADC graph, a *policy route* of any pair of PMs  $(pm_i, pm_j)$  is a path or walk starting  $pm_i$ , visiting all the  $n$  MBs at least once, and ending at  $pm_j$ . An OPR of  $(pm_i, pm_j)$  gives the minimum cost, denoted as  $opr(i, j)$ , among all its policy routes.  $\square$

OPR of  $(pm_i, pm_j)$  is essentially the shortest *s-t Hamiltonian path* [23] with  $s = pm_i$  and  $t = pm_j$  in a complete graph of  $pm_i, pm_j$  and all MBs (when  $pm_i = pm_j$ , it is a Hamiltonian cycle). The existing algorithm achieves approximation ratio of  $\frac{5}{3}$  and takes  $O(n^3)$  [23], where  $n$  is the number of MBs. Below we instead propose another  $O(n^3)$  but simpler algorithm to compute a policy route for  $(pm_i, pm_j)$  and show it has an approximation ratio of 2.

**Algorithm 2:** Compute A Policy Route for  $(pm_i, pm_j)$ .

**Input:** A PADC graph  $G$ , a pair of PMs  $(pm_i, pm_j)$ ,

and an unordered policy  $\{mb_1, mb_2, \dots, mb_n\}$ .

**Output:**  $pr(i, j)$ , cost of a policy route for  $(pm_i, pm_j)$ .

1.  $V_K^{i,j} = \{pm_i, pm_j, sw_1, sw_2, \dots, sw_n\}$ ;
2. Construct complete graph  $K^{i,j} = (V_K^{i,j}, E_K^{i,j})$ . For edge  $(u, v) \in E_K^{i,j}$ , its cost  $d(u, v)$  is the cost of the shortest path connecting  $u$  and  $v$  in  $G$ ;
3. Compute a minimum spanning tree MST for  $K^{i,j}$ ;
4. Compute a walk  $W$  from  $pm_i$  to  $pm_j$  on MST that visits all vertices in MST using each edge at most twice. Let the cost of  $W$  be  $pr(i, j)$ ;
5. **RETURN**  $pr(i, j)$ .

Using Algo. 2, Fig. 5 shows in blue dashed lines all three possible policy routes in the linear PADC of Fig. 3.

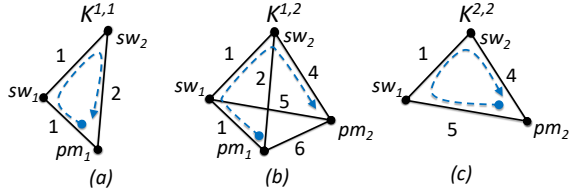


Fig. 5. How Algo. 3 works for PADC in Fig. 3. Blue dashed lines show the policy route of (a)  $(pm_1, pm_1)$ , (b)  $(pm_1, pm_2)$ , and (c)  $(pm_2, pm_2)$ .

**Lemma 1:**  $pr(i, j) \leq 2 \cdot opr(i, j)$ ,  $\forall pm_i, pm_j \in V_p$ .

**Proof:** Denote the cost of the MST computed in line 3 of Algo. 2 as  $c(\text{MST})$ ,  $c(\text{MST}) \leq opr(i, j)$ . Since the walk  $W$  found in line 4 uses each edge of the MST at most twice,  $pr(i, j) \leq 2 \cdot c(\text{MST})$ . Thus we have  $pr(i, j) \leq 2 \cdot opr(i, j)$ . ■

Next we present our PAL algorithm Algo. 3. It first computes the policy routes for all the  $|V_p| \cdot (|V_p| + 1) / 2$  pair of PMs using Algo. 2 and orders them in the non-descending order of their costs (lines 1-8). It then places the VM pairs (in the non-ascending order of their traffic rates) onto the first available PM pair, and updates the total communication cost accordingly (lines 9-21). Running time of Algo. 3 is  $O(|V_p|^2 \cdot n^3 + l)$ .

**Algorithm 3:** PAL Algorithm for Unordered Policy.

**Input:** A PADC with unordered policy  $\{mb_1, mb_2, \dots, mb_m\}$ ,

VM pairs  $\mathcal{P}$ ,  $V_p = \{pm_i\}$ , resource capacity  $rc_i$ .

**Output:** A placement  $p$  and the total comm. cost  $C_c(p, \vec{\pi})$ .

**Notations:**  $X$ : all PM pairs with their policy route costs.

$avail(pm_i)$ : available resource slots at  $pm_i$ , initially  $rc_i$ .

1.  $X = \emptyset$ ;  $avail(pm_i) = rc_i$ ,  $\forall pm_i \in V_p$ ;
2. **for**  $(i = 1; i \leq |V_p|; i++)$
3.   **for**  $(j = i; j \leq |V_p|; j++)$
4.     Compute  $pr(i, j)$  using Algo. 2;
5.      $X = X \cup \{(i, j, pr(i, j))\}$ ;
6.   **end for**;
7. **end for**;
8. Sort  $X$  in non-descending order of  $pr(i, j)$ . Let  $X$  be  $\{(s_1, t_1, pr(s_1, t_1)), (s_2, t_2, pr(s_2, t_2)), \dots\}$ ;
9.  $i = 1, j = 1, p = \emptyset, C_c(p, \vec{\pi}) = 0, \lambda_1 \geq \lambda_2 \dots \geq \lambda_l$ ;  
//  $i, j$ : indices for VM pairs and PM pairs respectively.
10. **while**  $(i \leq l)$  // not all VM pairs are placed yet
11.   **do**
12.     Place  $v_i$  at PM  $s_j$ ,  $v'_i$  at PM  $t_j$ ;
13.      $p = p \cup \{(s_j, t_j)\}$ ;
14.      $C_c(p, \vec{\pi}) += \lambda_i * pr(s_j, t_j)$ ;
15.      $avail(s_j) --, avail(t_j) --$ ;
16.      $i++$ ;
17.     **if**  $(i > l)$  **break**;
18.     **while**  $((s_j \neq t_j \wedge avail(s_j) \geq 1 \wedge avail(t_j) \geq 1)$   
       $\vee (s_j == t_j \wedge avail(s_j) \geq 2))$ ;
19.      $j++$ ; // the next available PM pairs
20.   **end while**;
21. **RETURN**  $p$  and  $C_c(p, \vec{\pi})$ .

**EXAMPLE 3:** Fig. 5 shows how Algo. 3 places  $(v_1, v'_1)$  and  $(v_2, v'_2)$ , with  $\vec{\lambda} = \langle 100, 1 \rangle$ , in the linear PADC in Fig. 3. It computes  $X = \{(1, 1, 4), (1, 2, 6), (2, 2, 10)\}$ . Thus  $(v_1, v'_1)$  is placed at  $pm_1$  and communicates in blue dashed line in Fig. 5(a). As  $pm_1$  is now full,  $(v_2, v'_2)$  is placed at  $pm_2$

and communicates in blue dashed line in Fig. 5(c). The total communication cost is  $100 \cdot 4 + 1 \cdot 10 = 410$ . □

**Theorem 3:** Algo. 3 achieves 2-approximation when  $l = 1$ .

**Proof:** Let the placement of  $(v_1, v'_1)$  computed by Algo. 3 be  $(pm_a, pm_b)$ . Let the optimal placement of  $(v_1, v'_1)$  be  $(pm_{a'}, pm_{b'})$  thus their optimal communication cost is  $opr(a', b')$ . From Lemma 1, we have  $pr(a', b') \leq 2 \cdot opr(a', b')$ . As the costs of all PM pair routes are sorted in non-descending order in Algo. 3,  $pr(a, b) \leq pr(a', b') \leq 2 \cdot opr(a', b')$ . ■

## V. PAM: POLICY-AWARE VM MIGRATION IN PADCS

### A. Ordered Policy.

1) *Problem Formulation:* In PAM, the initial VM placement is given by a *placement function*  $p : \mathcal{V} \rightarrow V_p$ , indicating that VM  $v \in \mathcal{V}$  is at PM  $p(v) \in V_p$ . The total communication cost of all the  $l$  VM pairs with placement  $p$  is thus  $C_c(p)$  (Eq. 1). Next, define a *VM migration function* as  $m : \mathcal{V} \rightarrow V_p$ , meaning that VM  $v$  will be migrated from PM  $p(v)$  to PM  $m(v)$  ( $m(v) = p(v)$  if  $v$  does not migrate). Let  $C_m(m) = \mu \cdot \sum_{i=1}^l (c(p(v_i), m(v_i)) + c(p(v'_i), m(v'_i)))$  be the *total migration cost* of all VM pairs with migration  $m$ . Let  $C_c(m)$  be the *total communication cost* of all VM pairs *after*  $m$ . Let  $C_t(m)$  be the *total migration and communication cost* of all pairs *after*  $m$ . Then  $C_t(m) = C_m(m) + C_c(m)$

$$\begin{aligned}
&= \sum_{i=1}^l \lambda_i \cdot \sum_{j=1}^{n-1} c(sw_j, sw_{j+1}) \\
&+ \sum_{i=1}^l \left( \mu \cdot c(p(v_i), m(v_i)) + \lambda_i \cdot c(m(v_i), sw_1) \right) \\
&+ \sum_{i=1}^l \left( \mu \cdot c(p(v'_i), m(v'_i)) + \lambda_i \cdot c(sw_n, m(v'_i)) \right). \quad (2)
\end{aligned}$$

The objective of PAM is to find a VM migration  $m$  that minimizes  $C_t(m)$  while satisfying resource constraint of PMs:  $|\{v \in \mathcal{V} | m(v) = pm_i\}| \leq rc_i, \forall pm_i \in V_p$ . As the first term on the right hand side in Eq. 2 is fixed under ordered policy, we only need to minimize the sum of the last two terms. Below we show that PAM under ordered policy is equivalent to minimum cost flow problem [5] in a properly transformed flow network.

2) *Minimum Cost Flow (MCF) Problem:* MCF is formally defined as follows. Let  $G = (V, E)$  be a directed graph. Denote the capacity and cost of an edge  $(u, v) \in E$  as  $ca(u, v)$  and  $d(u, v)$ , respectively. The amount of supply from source node  $s \in V$  equals the amount of demand at sink node  $t \in V$ . Denote a flow on edge  $(u, v)$  as  $f(u, v)$ ,  $f : E \rightarrow \mathbb{R}^+$ .  $f(u, v)$  is subject to (a) capacity constraint:  $f(u, v) \leq ca(u, v), \forall (u, v) \in E$  and (b) flow conservation constraint:  $\sum_{u \in V} f(u, v) = \sum_{u \in V} f(v, u)$ , for each  $v \in V \setminus \{s, t\}$ . The goal of MCF is to find a flow function  $f$  such that the total cost of the flow  $\sum_{(u, v) \in E} (d(u, v) \cdot f(u, v))$  is minimized. MCF can be solved efficiently and optimally by many combinatorial algorithms [5]. In this paper, we adopt the scaling push-relabel algorithm proposed by Goldberg [21] as it works well over a wide range of problem classes. The algorithm has the time complexity of  $O(A^2 \cdot B \cdot \log(A \cdot C))$ ,

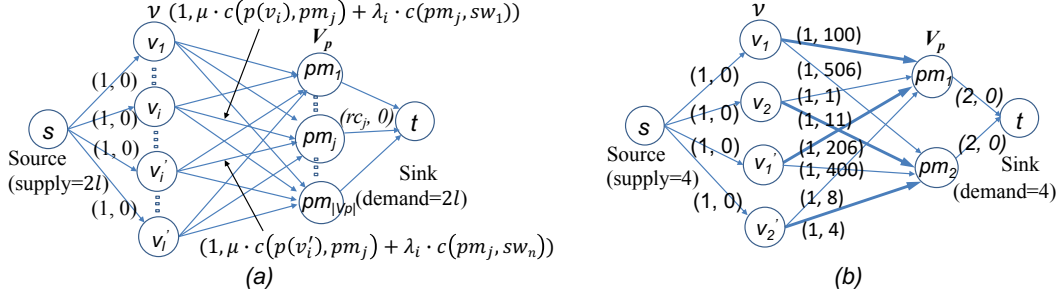


Fig. 6. (a) PAM under ordered policy is equivalent to MCF. (b) Graph transformation and MCF results (highlighted lines) for the PADC in Fig. 3(a).

where  $A$ ,  $B$ , and  $C$  are the number of nodes, number of edges, and maximum edge capacity in the flow network, respectively.

**Transforming a PADC to a Flow Network.** We transform a PADC  $G(V, E)$  into a flow network  $G'(V', E')$  following below five steps, as shown in Fig. 6(a).

Step I.  $V' = \{s\} \cup \{t\} \cup \mathcal{V} \cup V_p$ , where  $s$  is the source node and  $t$  is the sink node in the flow network.

Step II.  $E' = \{(s, v)\} \cup \{(v, pm_j)\} \cup \{(pm_j, t)\}$ , where  $v \in \mathcal{V}$  and  $pm_j \in V_p$ . Note that it is a complete bipartite graph between  $\mathcal{V}$  and  $V_p$ .

Step III. For each edge  $(s, v)$ , set its capacity as 1 and cost as 0. For each edge  $(pm_j, t)$ , set its capacity as  $rc_j$ , the resource capacity of  $pm_j$ , and cost as 0.

Step IV. For each edge  $(v_i, pm_j)$ , set its capacity as 1 and cost as  $\mu \cdot c(p(v_i), pm_j) + \lambda_i \cdot c(pm_j, sw_1)$ . For each edge  $(v'_i, pm_j)$ , set its capacity as 1 and cost as  $\mu \cdot c(p(v'_i), pm_j) + \lambda_i \cdot c(pm_j, sw_n)$ .

Step V. Set the supply at  $s$  and the demand at  $t$  as  $2l$ .

**Theorem 4:** PAM in ordered policy is equivalent to MCF in  $G'(V', E')$  thus can be solved optimally and efficiently.

**Proof Sketch:** Due to space constraint, we only give a high level sketch of the proof. By applying MCF algorithm upon the above flow network, it is able to achieve that a) every VM in the  $l$  VM pairs is assigned to exactly one resource slot in a PM while b) satisfying the resource capacity constraints of PMs and c) achieving the minimum total migration and communication cost for all the  $l$  VM pairs. ■

**Time Complexity.** As the number of nodes, edges, and maximum edge capacity in  $G'(V', E')$  are  $\bar{m} \cdot |V_p|$ ,  $\bar{m} \cdot |V_p|^2$ , and  $\bar{m}$  respectively, the MCF takes  $O(\bar{m}^3 \cdot |V_p|^4 \cdot \log(\bar{m}^2 \cdot |V_p|))$ .

**EXAMPLE 4:** Fig 6(b) illustrates how above transformation and MCF work for the same PADC in Fig. 3(a). MCF migrates  $v_1$  and  $v'_1$  to  $pm_1$ , and  $v_2$  and  $v'_2$  to  $pm_2$ , with total cost of  $100+11+206+4+101=422$ . Here 101 is the total communication cost between ingress switch  $sw_1$  and egress switch  $sw_2$ . This migration reduces the total cost of 606 before migration by 30%. Note that as  $v_1$  is initially located at  $pm_1$  and  $v'_2$  at  $pm_2$ , only  $v'_1$  and  $v_2$  actually migrate. □

3) *State-of-the-Art VM Migration Scheme:* Cui et al. [15] proposed a policy-aware VM management scheme named PLAN. The core concept of PLAN is *utility* of a VM migration. It is defined as a VM's communication cost reduction due to migration minus its migration cost. The goal of PLAN is to find a migration scheme that maximizes the *total utility*

of migrating all the VMs. PLAN is a greedy algorithm that works in rounds. In each round, it computes that which VM is migrated to which PM with available resources, such that the utility of this migration is the maximum among all the VMs that have not been migrated. This continues until all the VMs are migrated, or no more VM migration gives any positive utility. For the two VM pairs in Fig. 3(a), as there is no available resource slots, the migration cannot take place for PLAN. PLAN is however a heuristic algorithm that does not offer any performance guarantee. We prove in Lemma 2 below that its goal of maximizing the total utilities is equivalent to our goal of minimizing total communication and migration cost in PAM, thus we can compare our algorithms with PLAN.

**Lemma 2:** Minimizing total cost  $C_t(m)$  in PAM is equivalent to maximizing total utility in PLAN.

**Proof:** Denote the utility of migrating VM  $v$  as  $u(v)$ . Under migration function  $m$ , the *utility of migrating*  $v_i$  from its current PM  $p(v_i)$  to another PM  $m(v_i)$  is the reduction of its communication cost to the ingress switch minus the incurred migration cost. Thus  $u(v_i) = \lambda_i \cdot (c(p(v_i), sw_1) - c(m(v_i), sw_1)) - \mu \cdot c(p(v_i), m(v_i))$ . Similarly,  $u(v'_i) = \lambda_i \cdot (c(p(v'_i), sw_n) - c(m(v'_i), sw_n)) - \mu \cdot c(p(v'_i), m(v'_i))$ .

Given a  $p$  and a  $\vec{\lambda}$ , the total communication cost of the VMs  $C_c(p)$  can be computed using Eq. 1. Thus minimizing  $C_t(m)$  is equivalent to maximizing  $C_c(p) - C_t(m) \stackrel{\text{Eqs. 1,2}}{=} \sum_{i=1}^l \lambda_i \cdot (c(p(v_i), sw_1) + c(sw_n, p(v'_i)) - c(m(v_i), sw_1) - c(sw_n, m(v'_i))) - \mu \cdot \sum_{v \in \mathcal{V}} c(p(v), m(v)) = \sum_{i=1}^l (u(v_i) + u(v'_i))$ , which is the *total utility* of migrating all the VMs. ■

## B. Unordered Policy.

1) *Problem Formulation:* Under unordered policy, besides a VM migration function  $m : \mathcal{V} \rightarrow V_p$ , it defines for each VM pair  $(v_i, v'_i)$  an MB traversal function  $\pi^i : [1, 2, \dots, n] \rightarrow [1, 2, \dots, n]$ .  $\pi^i$  is a permutation function indicating that after VM migration, the  $j^{\text{th}}$  MB that  $(v_i, v'_i)$  visits is  $mb_{\pi^i(j)}$ . Let  $\vec{\pi} = \langle \pi^1, \pi^2, \dots, \pi^l \rangle$  and let  $C_t(m, \vec{\pi})$  denote the *total cost* of all the VM pairs with  $m$  and  $\vec{\pi}$ . Then  $C_t(m, \vec{\pi}) =$

$$\sum_{i=1}^l \left( \mu \cdot c(p(v_i), m(v_i)) + \mu \cdot c(p(v'_i), m(v'_i)) \right) + \sum_{i=1}^l \lambda_i \cdot \left( \sum_{j=1}^{n-1} c(sw(\pi^i(j)), sw(\pi^i(j+1))) \right) + c(m(v_i), sw(\pi^i(1))) + c(sw(\pi^i(n)), m(v'_i)). \quad (3)$$

The first and second terms in Eq. 3 are the total migration cost and total communication cost respectively. The objective of PAM under unordered policy is to find an  $m$  and a  $\vec{\pi}$  to minimize  $C_t(m, \vec{\pi})$  while satisfying resource constraints of PMs:  $|\{v \in \mathcal{V} | m(v) = i\}| \leq rc_i, \forall i \in V_p$ .

2) *VM Migration Algorithm for Unordered Policy*: Algo. 4 below first computes costs for all the  $|V_p| \cdot (|V_p| + 1)/2$  policy routes (lines 2-6). Then for each VM pair (in the non-ascending order of their traffic rates), it finds a PM pair to migrate to, such that the resulted cost for this VM pair is the minimum among all the unassigned VM pairs in this round (lines 7-22). After the entire migration scheme  $m$  is computed, it finally migrates the VMs and returns the total cost (lines 23 and 24). Its takes  $O(|V_p|^2 \cdot (n^3 + l))$ .

**Algorithm 4:** PAM Algorithm for Unordered Policy.

**Input:** A PADC with unordered policy  $\{mb_1, mb_2, \dots, mb_n\}$ ,  $V_p = \{pm_i\}$ , resource capacity  $rc_i$ , VM pair placement  $p$ .

**Output:** A migration scheme  $m$  and the total cost  $C_t(m, \vec{\pi})$ .

**Notations:**  $i, j$ : indices for PM pairs;  $k$ : index for VM pairs.

$c_{i,j}$ : the total cost of a VM pair if its source VM is migrated to  $pm(i)$  and destination VM to  $pm(j)$ .

$a, b$ : indices of a PM pair that gives minimum total cost.

$avail(pm_i)$ : number of available slots at  $pm_i$ , initially  $rc_i$ .

1.  $m = \phi, C_t(m, \vec{\pi}) = 0, k = 1, \lambda_1 \geq \lambda_2 \dots \geq \lambda_l$ ;
2. **for** ( $i = 1; i \leq |V_p|; i++$ )
3.   **for** ( $j = i; j \leq |V_p|; j++$ )
4.     Compute  $pr(i, j)$  using Algo. 2;
5.   **end for**;
6. **end for**;
7. **while** ( $k \leq l$ )   // find PM pair for VM pair  $(v_k, v'_k)$
8.    $c_{min} = \infty$ ;   // minimum total cost for  $(v_k, v'_k)$
9.   **for** ( $i = 1; i \leq |V_p|; i++$ )
10.     **if** ( $avail(pm_i) == 0$ )   **continue**; //  $pm_i$  is full
11.     **for** ( $j = i; j \leq |V_p|; j++$ )
12.       **if** ( $(avail(pm_j) == 0) \vee$    //  $pm_j$  is full  
        $(i == j \wedge avail(pm_j) \leq 1)$ ) **continue**;
13.        $c_{i,j} = 0$ ;
14.        $c(pm_i) = \mu \cdot c(p(v_k), pm_i)$ , // cost of migrating  $v_k$   
        $c(pm_j) = \mu \cdot c(p(v'_k), pm_j)$ ; // to  $pm_i, v'_k$  to  $pm_j$
15.        $c_{i,j} = \lambda_k \cdot pr(i, j) + c(pm_i) + c(pm_j)$ ;
16.       **if** ( $c_{i,j} < c_{min}$ )    $a = i, b = j, c_{min} = c_{i,j}$ ;
17.     **end for**;
18.   **end for**;
19.    $m = m \cup \{(pm_a, pm_b)\}$ ; // update migration scheme  $m$
20.    $C_t(m, \vec{\pi}) += c_{min}$ ;   // update total cost
21.    $avail(pm_a) --, avail(pm_b) --$ ;
22. **end while**;
23. Migrate  $(v_1, v'_1), \dots, (v_l, v'_l)$  according to  $m$ ;
24. **RETURN**  $m$  and  $C_t(m, \vec{\pi})$ .

**EXAMPLE 5:** For the two VM pairs stored in the PADC of Fig. 3 (a), Algo. 4 will migrate both  $v_1$  and  $v'_1$  to  $pm_1$ , resulting in cost of 406 for this pair. As  $pm_1$  is now full, it will then migrate both  $v_2$  and  $v'_2$  to  $pm_2$ , resulting in cost of 16 for this pair. The total cost of the two pairs is 422.  $\square$

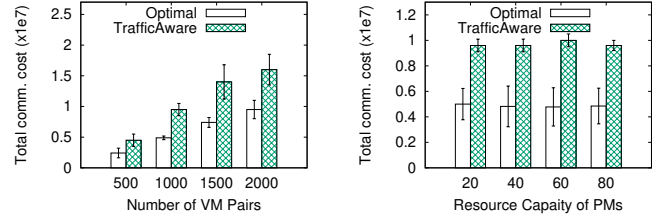
**Theorem 5:** Under unordered policy, PAM is NP-hard even

for one pair of VMs (i.e.,  $l = 1$ ).

**Proof:** We reduce a variation of the *s-t traveling salesman path problem* to this special case. By variation, we mean that in complete graph  $K = (V_K, E_K)$ , each node in  $V_K$  has a cost. Thus the cost of the *s-t* shortest Hamiltonian path includes the costs of  $s$  and  $t$ . The rest of the proof is then similar to that in Theorem 2 with one augmentation: For  $pm_i$ , its cost  $c(pm_i)$  is the migration cost of  $v_1$  from  $p(v_1)$  to  $pm_i$ ; for  $pm_j$ , its cost  $c(pm_j)$  is the migration cost of  $v'_1$  from  $p(v'_1)$  to  $pm_j$ .  $\blacksquare$

**Theorem 6:** Algo. 4 achieves 2-approximation when  $l = 1$ .

**Proof:** For Algo. 4, let the PM pair that  $(v_1, v'_1)$  migrate to be  $(pm_a, pm_b)$ . Let their optimal VM migration be  $(pm_{a'}, pm_{b'})$  and their optimal total cost be  $C_t^{opt}(m, \vec{\pi})$ . The total cost of  $(v_1, v'_1)$  computed by Algo. 4  $C_t(m, \vec{\pi}) = \lambda_1 \cdot pr(a, b) + \mu \cdot c(p(v_1), pm_a) + \mu \cdot c(p(v'_1), pm_b) \leq \lambda_1 \cdot pr(a', b') + \mu \cdot c(p(v_1), pm_{a'}) + \mu \cdot c(p(v'_1), pm_{b'}) \leq 2 \cdot \lambda_1 \cdot opr(a', b') + 2 \cdot \mu \cdot c(p(v_1), pm_{a'}) + 2 \cdot \mu \cdot c(p(v'_1), pm_{b'}) = 2 \cdot C_t^{opt}(m, \vec{\pi})$ .  $\blacksquare$



(a) Varying  $l$ .  $n = 3, rc = 40$ .

(b) Varying  $rc$ .  $l = 1000, n = 3$ .

Fig. 7. Comparing with TrafficAware in ordered policy,  $k = 16$ .

## VI. PERFORMANCE EVALUATION

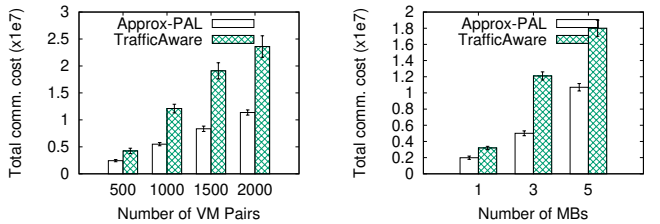
In this section we compare our algorithms with existing work (Table II). For PAL, we name the optimal algorithm for ordered policies (Algo. 1) as **Optimal** and the approximation algorithm for unordered (Algo. 3) as **Approx-PAL**, and compare them with **TrafficAware** [33], a seminal but policy-oblivious VM placement algorithm. For PAM, we refer to the minimum cost flow-based algorithm for ordered policy as **MCF** and the approximation algorithm for unordered (Algo. 4) as **Approx-PAM**, and compare them with **PLAN** [15].

TABLE II  
COMPARING PAM AND PAL ALGORITHMS.

	Ordered Policy	Unordered Policy	Existing Work
PAL	Optimal	Approx-PAL	TrafficAware [33]
PAM	MCF	Approx-PAM	PLAN [15]

We consider fat-tree PADCs of size  $k = 8$  with 128 PMs and size  $k = 16$  with 1024 PMs. The traffic rates of VM pairs are in the range of  $[0, 1000]$  – Following flow characteristics found in Facebook data centers (Section 5.1, [36]), 25% of VM pairs have light traffic rates in  $[0, 300]$ , 70% medium traffic rates in  $[300, 700]$ , and 5% heavy rates in  $(700, 1000]$ . As suggested by Cisco design guide [3], we install a number of MBs on different aggregation switches in the PADC. As 80% of cloud data center traffic originated by servers stays within the rack [9], for the initial VM placement in PAM, we place 80% of the VM pairs into the PMs under the same edge switches while the rest 20% under different edge switches. Each data point in the plots is an average of 20 runs with





(a) Varying  $l$ .  $n = 3$ ,  $rc = 40$ . (b) Varying  $n$ .  $l = 1000$ ,  $rc = 40$ .  
Fig. 8. Comparing with TrafficAware in unordered policy,  $k = 16$ .

95% confidence interval. In each run a new set of VM pairs are placed (for PAL) or migrated (for PAM) in the PADC.

**Comparing with TrafficAware.** As TrafficAware only assigns VMs to the same PMs or PMs in close proximity, and does not consider the proximity of the PMs to the MBs, we implement TrafficAware as follows for fair comparison. In ordered policy, it places VM pairs (in non-ascending order of their traffic rates) to the PMs that are closest to the ingress switch. In unordered policy, it works like Algo. 3 but only considers the Hamiltonian cycle case, as TrafficAware always places VM pairs in the same PM if possible. For ordered policy, Fig. 7(a) varies the number of VM pairs  $l$  and shows that Optimal yields 46-49% less costs than TrafficAware. Fig. 7(b) varies resource capacities of PMs  $rc$  and shows that Optimal outperforms TrafficAware by around 48%. Fig. 8 compares Approx-PAL and TrafficAware under unordered policy. It varies  $l$  as well as number of MBs  $n$  and shows that Approx-PAL outperforms TrafficAware by 37-58% in all scenarios. Above results are evident as Optimal is optimal and Approx-PAL is 2-approximation policy-aware algorithms while TrafficAware is policy-oblivious, inducing enormous traffic when VM communication traverses the MBs.

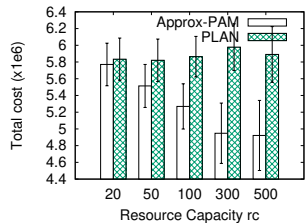
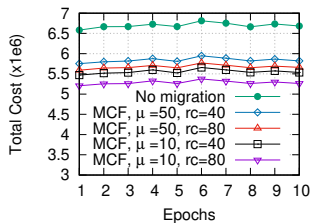
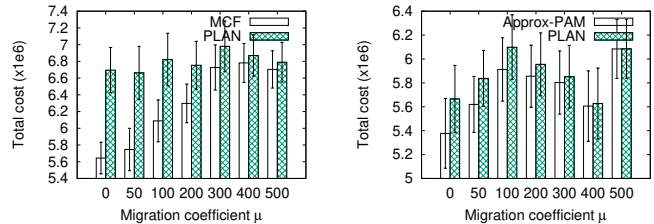


Fig. 9. Effects of VM migration. Fig. 10. Comparing with PLAN.

**Effects of VM Migrations.** Fig. 9 investigates how much cost reduction VM migration brings to a PADC ( $k = 8$ ,  $l = 1000$ ,  $n = 3$ ) compared to without migration. VM migrations take place in *epochs*. At the beginning of each epoch, VM pairs change their traffic rates to new values in  $[0,1000]$  following aforesaid Facebook flow pattern. For migrations, it then executes MCF and calculates the total migration and communication cost. For no migration, it simply recalculates the total communication cost using the new traffic rates. We set the migration coefficient  $\mu$  as 10 and 50, and let the PADC run continuously for ten epochs. Fig. 9 shows the total cost of VM pairs in each epoch with and without VM migration, for both  $rc = 40$  and 80. The cost of  $\mu = 50$  is larger than that of  $\mu = 10$ , as migration costs increase



(a) Ordered Policy. (b) Unordered Policy.  
Fig. 11. Comparing with PLAN,  $k = 8$ ,  $l = 1000$ ,  $n = 3$ ,  $rc = 40$ .

with the increase of  $\mu$ . In either case, the cost of  $rc = 80$  is smaller than that of  $rc = 40$ , as there are more resource slots available to achieve cost-efficient VM migrations. In all cases, VM migration reduces the total costs by up to 25% ( $\mu = 10$  and  $rc = 80$ ) compared to without migration.

**Comparing with PLAN.** We then compare our PAM algorithms with PLAN by increasing migration coefficient  $\mu$ . As PLAN is only designed for ordered policy, for the purpose of comparison in unordered policy, we implement it as below greedy algorithm. For each VM pair, it finds the MB closest to the source VM as ingress MB and the one closest to the destination VM as egress MB. It then finds a MB sequence by starting from the ingress MB, visiting an unvisited closest MB, so on and so forth until all the MBs are visited, and finally visiting the egress MB. Fig. 11(a) shows that under ordered policy, the MCF outperforms the PLAN by around 20% when  $\mu$  is small. With the increase of  $\mu$ , PLAN and MCF start to perform close due to high migration cost. Fig. 11(b) shows that under unordered policy, Approx-PAM outperforms PLAN slightly for the entire range of  $\mu$ . Finally, Fig. 10 compares Approx-PAM and PLAN by varying  $rc$  while fixing  $\mu$  as 20, and shows that Approx-PAM outperforms PLAN for the entire range of  $rc$ . With the increase of the  $rc$ , the performance difference between Approx-PAM and Greedy gets larger, around 30%.

## VII. CONCLUSIONS AND FUTURE WORK

We proposed and studied PAL and PAM: VM placement and VM migration in the PADC, an emerging cloud computing platform. We demonstrated that VM migration is an effective technique to alleviate dynamic VM traffic in PADCs. Working together, PAL and PAM place and then migrate VMs in the event of dynamic traffic fluctuation, achieving optimal and near-optimal network resource management for a PADC's lifetime. We uncovered a suite of new policy-aware problems and designed optimal, approximation, and heuristic algorithms. We also showed that our results outperform the state-of-the-arts. As future work, we will study if the optimality and approximability of our algorithms still hold when VMs have different resource demands. We will also study how VNF migration mitigates diverse and dynamic traffic and design a holistic VNF+VM migration scheme to achieve ultimate resource optimization in PADCs.

## ACKNOWLEDGMENT

This work was supported by NSF Grant CNS-1911191.

## REFERENCES

- [1] Amazon lex. <https://aws.amazon.com/lex/>.
- [2] Cisco global cloud index: Forecast and methodology, 2016 to 2021 white paper. <https://www.cisco.com/c/en/us/solutions/service-provider/global-cloud-index-gci/white-paper-listing.html>.
- [3] Cisco virtualized multi-tenant data center, version 2.0 compact pod design guide. <http://hyperurl.co/hpj2xt>.
- [4] Slack, the collaboration software that moves work forward. <http://slack.com>.
- [5] R. K. Ahuja, T. L. Magnanti, and J. B. Orlin. *Network Flows: Theory, Algorithms, and Applications*. Prentice-Hall, Inc., 1993.
- [6] M. Al-Fares, A. Loukissas, and A. Vahdat. A scalable, commodity data center network architecture. *SIGCOMM Comput. Commun. Rev.*, 38(4):63–74, 2008.
- [7] M. Alicherry and T.V. Lakshman. Optimizing data access latencies in cloud systems by intelligent virtual machine placement. In *Proc. of IEEE INFOCOM 2013*.
- [8] A. Beloglazov and R. Buyya. Managing overloaded hosts for dynamic consolidation of virtual machines in cloud data centers under quality of service constraints. *IEEE Transactions on Parallel and Distributed Systems*, 24(7):1366 – 1379, 2013.
- [9] T. Benson, A. Akella, and D. A. Maltz. Network traffic characteristics of data centers in the wild. In *Proc. of ACM IMC 2010*.
- [10] B. Carpenter and S. Brim. Middleboxes: Taxonomy and issues, 2002. <https://tools.ietf.org/html/rfc3234>.
- [11] C. Clark, K. Fraser, and S. Hand. Live migration of virtual machines. In *Proc. of NSDI 2005*.
- [12] R. Cohen, L. Lewin-Eytan, J. Seffi Naor, and D. Raz. Almost optimal virtual machine placement for traffic intense data centers. In *Proc. of IEEE INFOCOM 2013*.
- [13] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. The MIT Press, 2009.
- [14] L. Cui, R. Cziva, F. P. Tso, and D. P. Pezaros. Synergistic policy and virtual machine consolidation in cloud data centers. In *Proc. of IEEE INFOCOM 2016*.
- [15] L. Cui, F. P. Tso, D. P. Pezaros, W. Jia, and W. Zhao. Plan: Joint policy- and network-aware vm management for cloud data centers. *IEEE Transactions on Parallel and Distributed Systems*, 28(4):1163–1175, 2017.
- [16] Y. Cui, Z. Yang, S. Xiao, X. Wang, and S. Yan. Traffic-aware virtual machine migration in topology-adaptive dcn. *IEEE/ACM Transactions on Networking*, 25(6):3427 – 3440, 2017.
- [17] V. Eramo, E. Miucci, M. Ammar, and F. G. Lavacca. An approach for service function chain routing and virtual function network instance migration in network function virtualization architectures. *IEEE/ACM Transactions on Networking*, 25(4):2008–2025, 2017.
- [18] J. Fan, C. Guan, Y. Zhao, and C. Qiao. Availability-aware mapping of service function chains. In *Proc. of IEEE INFOCOM 2017*.
- [19] S. K. Fayazbakhsh, L. Chiang, V. Sekar, M. Yu, and J. C. Mogul. Enforcing network-wide policies in the presence of dynamic middlebox actions using flowtags. In *Proc. of USENIX NSDI 2014*.
- [20] A. Gember-Jacobson, R. Viswanathan, C. Prakash, R. Grandl, J. Khalid, S. Das, and A. Akella. Opennf: Enabling innovation in network function control. In *Proc. of ACM SIGCOMM 2014*.
- [21] A. V. Goldberg. An efficient implementation of a scaling minimum-cost flow algorithm. *J. Algorithms*, 22:1–29, 1997.
- [22] A. Gushchin, A. Walid, and A. Tang. Scalable routing in sdn-enabled networks with consolidated middleboxes. In *Proc. of ACM Hotmiddlebox*, 2015.
- [23] J.A. Hoogeveen. Analysis of christofides’ heuristic: Some paths are more difficult than cycles. *Operations Research Letters*, 10:291 – 295, 1991.
- [24] H. Huang, S. Guo, J. Wu, and J. Li. Service chaining for hybrid network function. *IEEE Transactions on Cloud Computing*, 7:1082–1094, 2019.
- [25] Y. Jiang, Y. Cui, W. Wu, Z. Xu, J. Gu, K. K. Ramakrishnan, Y. He, and X. Qian. Speedybox: Low-latency nfv service chains with cross-nf runtime consolidation. In *Proc. of IEEE ICDCS 2019*.
- [26] D. A. Joseph, A. Tavakoli, and I. Stoica. A policy-aware switching layer for data centers. In *Proc. of ACM SIGCOMM 2008*.
- [27] T. Kuo, B. Liou, K. C. Lin, and M. Tsai. Deploying chains of virtual network functions: On the relation between link and server usage. *IEEE/ACM Transactions on Networking*, 26(4):1562–1576, Aug 2018.
- [28] L. E. Li, V. Liaghat, H. Zhao, M. Hajiaghayi, D. Li, G. Wilfong, Y. R. Yang, and C. Guo. Pace: Policy-aware application cloud embedding. In *Proc. of IEEE INFOCOM 2013*.
- [29] X. Li, J. Wu, S. Tang, and S. Lu. Let’s stay together: Towards traffic aware virtual machine placement in data centers. In *Proc. of IEEE INFOCOM 2014*.
- [30] J. Liu, Y. Li, Y. Zhang, L. Su, and D. Jin. Improve service chaining performance with optimized middlebox placement. *IEEE Transactions on Services Computing*, 10(4):560–573, 2017.
- [31] J. Liu, W. Lu, F. Zhou, P. Lu, and Z. Zhu. On dynamic service function chain deployment and readjustment. *IEEE Transactions on Network and Service Management*, 14(3):543–553, 2017.
- [32] V. Mann, A. Gupta, P. Dutta, A. Vishnoi, P. Bhattacharya, R. Poddar, and A. Iyer. Remedy: Network-aware steady state vm management for data centers. In *Proc. of the NETWORKING 2012*.
- [33] X. Meng, V. Pappas, and L. Zhang. Improving the scalability of data center networks with traffic-aware virtual machine placement. In *Proc. of IEEE INFOCOM 2010*.
- [34] A. Mohammadkhan, S. Ghapani, G. Liu, W. Zhang, K. K. Ramakrishnan, and T. Wood. Virtual function placement and traffic steering in flexible and dynamic software defined networks. In *The 2015 IEEE International Workshop on Local and Metropolitan Area Networks*.
- [35] Z. A. Qazi, C.-C. Tu, L. Chiang, R. Miao, V. Sekar, and M. Yu. Simplefying middlebox policy enforcement using sdn. In *Proc. of ACM SIGCOMM 2013*.
- [36] A. Roy, H. Zeng, J. Bagga, G. Porter, and A. C. Snoeren. Inside the social network’s (datacenter) network. In *Proc. of ACM SIGCOMM 2015*.
- [37] G. Sallam, G.R. Gupta, B. Li, and B. Ji. Shortest path and maximum flow problems under service function chaining constraints. In *Proc. of IEEE INFOCOM 2018*.
- [38] Y. Sang, B. Ji, G. R. Gupta, X. Du, and L. Ye. Provably efficient algorithms for joint placement and allocation of virtual network functions. In *Proc. of Infocom 2017*.
- [39] V. Sekar, N. Egi, S. Ratnasamy, M. K. Reiter, and G. Shi. Design and implementation of a consolidated middlebox architecture. In *Proc. of NSDI 2012*.
- [40] J. Sherry, S. Hasan, C. Scott, A. Krishnamurthy, S. Ratnasamy, and V. Sekar. Making middleboxes someone else’s problem: Network processing as a cloud service. In *Proc. of ACM SIGCOMM 2012*.
- [41] V. Shrivastava, P. Zerkos, K.-W. Lee, H. Jamjoom, Y.H. Liu, and S. Banerjee. Application-aware virtual machine migration in data centers. In *Proc. of INFOCOM 2011, mini conference*.
- [42] H. Wang, Y. Li, Y. Zhang, and D. Jin. Virtual machine migration planning in software-defined networks. In *Proc. of Infocom 2015*.
- [43] L. Yu, L. Chen, Z. Cai, H. Shen, Y. Liang, and Y. Pan. Stochastic load balancing for virtual resource management in datacenters. *IEEE Transactions on Cloud Computing (Early Access)*, 2018.
- [44] J. Zhang, F. Ren, and C. Lin. Delay guaranteed live migration of virtual machines. In *Proc. of INFOCOM 2014*.
- [45] Y. Zhang, N. Beheshti, L. Beliveau, G. Lefebvre, R. Manghirmalani, R. Mishra, R. Patney, M. Shirazipour, R. Subrahmaniam, C. Truchan, and M. Tatipamula. Steering: A software-defined networking for inline service chaining. In *Proc. of IEEE ICNP 2013*.