

# Kraken: Online and Elastic Resource Reservations for Cloud Datacenters

Carlo Fuerst, Stefan Schmid<sup>1</sup>, Lalith Suresh, and Paolo Costa

**Abstract**—In cloud environments, the absence of strict network performance guarantees leads to unpredictable job execution times. To address this issue, recently, there have been several proposals on how to provide guaranteed network performance. These proposals, however, rely on computing resource reservation schedules *a priori*. Unfortunately, this is not practical in today’s cloud environments, where application demands are inherently unpredictable, e.g., due to differences in the input data sets or phenomena, such as failures and stragglers. To overcome these limitations, we designed KRAKEN, a system that allows to dynamically update minimum guarantees for both network bandwidth and compute resources *at runtime*. Unlike previous work, Kraken does not require prior knowledge about the resource needs of the applications but allows to modify reservations at runtime. Kraken achieves this through an *online resource reservation scheme*, which comes with provable optimality guarantees. In this paper, we motivate the need for dynamic resource reservation schemes, present how this is provided by Kraken, and evaluate Kraken via extensive simulations and a preliminary Hadoop prototype.

**Index Terms**—Network virtualization, embedding, predictable performance, algorithms.

## I. INTRODUCTION

CLOUD-BASED applications, including batch processing, streaming, and scale-out databases, generate a significant amount of network traffic and a considerable fraction of their runtime is due to network activity. For example, traces of jobs from a Facebook cluster reveal that network transfers on average account for 33% of the execution time [24].

Unfortunately, as reported in previous studies [5], in existing cloud infrastructures the bandwidth available to the tenants varies significantly over time, i.e., by a factor of five or more [35], even within the same day. Given the time spent in network activity by these applications, this variability has a non-negligible impact on the application performance, which makes it impossible to accurately estimate the execution time in advance [26].

Manuscript received January 18, 2017; revised August 5, 2017 and October 27, 2017; accepted November 30, 2017; approved by IEEE/ACM TRANSACTIONS ON NETWORKING Editor X. Liu. Date of publication December 22, 2017; date of current version February 14, 2018. This work was supported in part by the Aalborg University’s Talent Management Programme through the Project PreLytics and in part by the German Ministry for Education and Research through the Berlin Big Data Center. (Corresponding author: Stefan Schmid.)

C. Fuerst was with the Technical University of Berlin, 10623 Berlin, Germany. He is now with SAP, 14469 Potsdam, Germany.

S. Schmid is with the University of Vienna, 1090 Vienna, Austria, and also with the Faculty of Engineering and Science, Department of Computer Science, Aalborg University, 9220 Aalborg, Denmark (e-mail: schmiste@gmail.com).

L. Suresh is with VMware Research, Palo Alto, CA 94304 USA.

P. Costa is with Microsoft Research, Cambridge CB1 2FB, U.K.

Digital Object Identifier 10.1109/TNET.2017.2782006

Over the last years, several solutions have been proposed to improve the sharing of network bandwidth among tenants,

by leveraging admission control and bandwidth reservations, thus enabling tenants to specify *absolute* guarantees [5], [9], [17], [21], [30], [31], [33]. In particular, many of these proposals offer a *virtual cluster* abstraction [5], [9], which provides the tenants with the illusion of having their own dedicated network. A virtual cluster *guarantees* a specified *minimal* bandwidth between all tenant’s virtual machines, independently of their locations in the datacenter topology.

However, the vast majority of existing solutions providing absolute bandwidth guarantees are based on *offline* and *constant* reservations schemes [5], [9], [17], [21], [28], [33]: they require that tenants announce the entire resource reservation schedule *ahead of time*, i.e., at job submission time. They typically assume that the corresponding resource reservations need to be *constant* over time, and hence tenants or operators either have to over-provision during idle times (thus reducing efficiency and inflating cost) or under-provision during peak times (thus reducing application performance), or both. Notable exceptions are Cicada [22], which offers predictive instead of absolute guarantees, and Proteus [9], which allows tenants to specify time-varying bandwidth reservations. However, even with Proteus, the reservations must be made at the startup time and they cannot be changed afterwards. This inflexibility is at odds with the cloud computing paradigm, which enables elasticity by allowing to “scale out” or “scale in” applications at runtime. We argue that in most cases it is very hard to accurately estimate application resource needs ahead of time, rendering offline reservation schemes inadequate. Several factors contribute to this unpredictability including unexpected events such as stragglers and failures [3], [4] as well as spikes in application demand (flash crowds). As a consequence, tenants would be likely to either over-estimate the network usage, which would consequently inflate their expenditure, or to under-estimate it, which would result in poor application performance, and, as previously observed [5], in potentially longer job running times, which would also lead to higher customer costs. Over-estimation of network resource would also negatively impact the cloud provider as it would reduce its ability to accommodate more clients, thus reducing its market share. In contrast, we believe that a more principled solution would be to enable re-configuring network reservations at runtime so as to better match tenant requirements and avoiding resource over- or under-provision.

A naive approach to enable runtime reconfiguration would be to restart the resource allocation from scratch every time

an update request is received. This, however, would introduce an unacceptable overhead as most (if not all) the compute resources such as VMs need be migrated. At the other extreme, there are approaches such as Blender [33] that support a weak form of reconfiguration by allowing tenants to update rate limiters at runtime. This, however, prevents users from upgrading both compute and network resources at the same time. More importantly, as we show in the evaluation, since no migration is considered, the efficacy of the solution is very limited. In this paper, we strike a balance between these two approaches by allowing users to dynamically reconfigure both compute and network resources *simultaneously* while minimizing the number of migrations.

### A. Our Contribution

We make the following contributions.

- 1) *The need for online resource reservation schemes:* We show that offline resource reservation schemes are insufficient: Even for simple Hadoop jobs, small internal changes can lead to significantly different executions. Therefore, in order to meet application performance goals, not only strict resource isolation needs to be provided, but also a possibility to update these reservations at runtime.
- 2) *The Kraken system:* We design Kraken, a system which supports the online (and joint) update of both bandwidth as well as the compute resources. Kraken can also perform migrations in order to satisfy upgrade requests: While the migration of entire virtual machines may be expensive in practice, Kraken only assumes that compute units, the endpoints of traffic flows, can be migrated. Kraken comes with provable performance guarantees and ensures (i) the satisfaction of all upgrade/downgrade requests for which this is possible, (ii) minimal reconfiguration and resource costs, (iii) low runtimes.
- 3) *Benefits of online resource reservations:* Our simulations show the benefits of elastic resource reservations. We also demonstrate the feasibility of our approach in practice, through a preliminary implementation of Kraken on top of Hadoop.

Kraken can be used for many applications that benefit from resource elasticity, including batch-processing applications (e.g., graph processing or distributed databases) or high-performance computing applications.

### B. Scope and Non-Goals

We focus on *how* to efficiently embed and reconfigure virtual clusters; a detailed discussion of *when* to change a virtual network specification is left for future work. The time and extent of upgrades and downgrades depend on the setting, on the type of application, as well as on the tenants' resp. operators' objectives. We also note on this occasion that Kraken is not online in the sense of competitive analysis [10].

In general, Kraken is agnostic to where the update requests come from: they can come from the operator itself (e.g., by monitoring network traffic or other application metrics similar to Amazon Auto scaling), from the application framework

(e.g., a Kraken-aware version of Hadoop could issue an update request when the shuffle phase is about to start/end) or could be triggered by the tenants themselves to improve application performance by removing network bottlenecks, in a way similar to how today's tenants can increase the number VMs rented to speed-up their applications or services. Kraken is not tied to any of these approaches and can support all of them. While this paper only focuses on providing the *mechanisms*, we think it is an interesting direction for future work to devise the *policies* that can be built on top of Kraken.

### C. Organization

We first motivate the online approach pursued in this paper in Section II. In Section III, we describe our model and give an example illustrating the challenge of dynamic resource reservation schemes. The Kraken system and its embedding and reconfiguration algorithms are presented in detail and analyzed formally in Section IV. We present our simulation results in Section V and report on a preliminary prototype in Hadoop in Section VI. After reviewing related work in Section VII, we conclude in Section VIII.

## II. MOTIVATION FOR AN ONLINE APPROACH

Before presenting our solution in detail, we argue that today's offline reservation schemes are not sufficient to ensure application performance guarantees in an efficient manner.

We distinguish between two offline reservation schemes: (1) schemes with *constant* resource reservations such as the ones proposed in [5] and [17]; and (2) schemes such as Proteus [9] with time-varying resource reservations which, however, need to be announced ahead of time and, hence, require accurately predicting a job's resource-utilization over time, e.g., using data from previous runs.

Constant reservation schemes are wasteful for any application with time-varying resource demands, such as MapReduce applications, which cycle between network-intensive and compute-intensive phases [9], or an online computer game whose demand is subject to time-of-day effects [35].

While offline and time-varying reservations may be possible in idealized conditions, in practice, this is rarely the case. This is obvious for continuously running applications, such as a web-service or video-on-demand service, whose popularity can change significantly and unexpectedly. But, as we show next, even the resource pattern of very simple MapReduce applications are hard to predict accurately. It has been reported that stragglers can be several times slower than the median task completion time [3], [4], [12], [20], [34]. Stragglers occur due to a variety of environmental factors such as slow disks and failures. Cluster frameworks typically use control loops based on these factors to (re-)schedule tasks, e.g., Hadoop's speculative executor. This makes it hard to predict if there will be stragglers in the first place and if so, when and where the cluster framework will re-schedule a slow task.

To highlight this, we rely on a very simple experiment wherein we run a Hadoop cluster in an OpenStack-based testbed. For this we use five physical servers (8 CPU cores and 64GB of RAM) with one virtual machine each. Each virtual

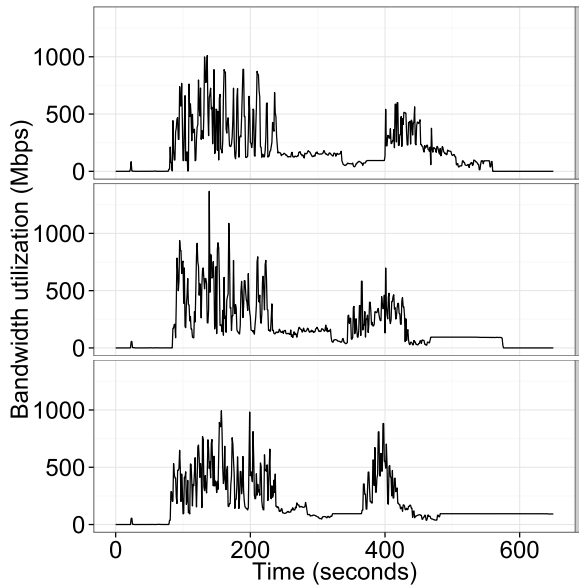


Fig. 1. Execution unpredictability—Bandwidth utilization for the same TeraSort workload with three different datasets of the same size generated via TeraGen. Bandwidth utilization over time varies between the runs, necessitating the need for online bandwidth reservation schemes as opposed to offline ones.

machine is allocated 4 virtual cores with 4 GB of RAM. Each node of the Hadoop cluster is mapped to a virtual machine each (one master, four slaves). The Hadoop workload consists of a TeraSort job, operating on 150 million 100-byte records. We repeat the experiment five times with speculative execution enabled. Figure 2 (left) indicates the variance in job completion times across the runs: a range of 150 seconds. This observation is also supported by Figure 2 (right) which indicates the number of straggling tasks that were speculatively re-executed by the Hadoop cluster.

Figure 1 indicates the bandwidth consumption in the cluster across three different runs of the TeraSort job executed without speculative execution enabled. In each run, the dataset is generated afresh using the TeraGen command. Note that the bandwidth utilization has varying profiles over time in each case. These observations serve to demonstrate that even with the same workload and a dataset of the same size being re-executed, it is difficult to predict how a job progresses over time.

Note that since TeraSort is IO-bound and all data are randomly generated with a uniform distribution, its behavior is much more regular than most other jobs used in data analytics, which can suffer from skewed data distribution, irregular computation patterns, etc. Therefore, we expect real jobs to exhibit even higher variance across runs, as it is often reported in literature [3], [4], [12], [20].

In conclusion, we argue that offline approaches for resource reservations such as Proteus do not suffice, as cloud environments such as Amazon EC2 [34] are likely to be much more noisy than our environment studied here. This makes it more difficult to predict performance of an application a priori, which underlines the need for dynamic and online reservation schemes. The need for online solutions is exacerbated in

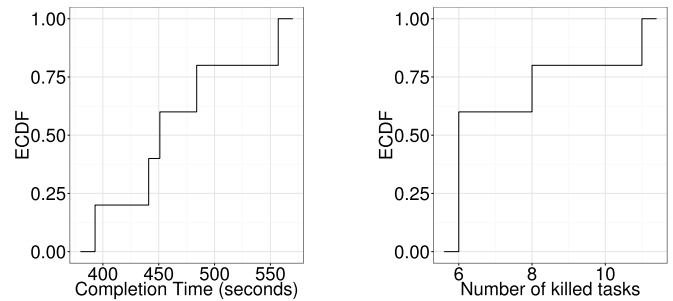


Fig. 2. Execution unpredictability—Completion times of jobs in the presence of speculative execution (left) and the number of speculated tasks (right).

systems in which demand can vary over time, e.g., long-running applications or streaming applications.

### III. MODEL & EXAMPLE

We start by introducing the settings and the virtual network abstraction considered in this paper, and subsequently highlight the algorithmic challenge.

#### A. Setting

We consider the standard *Virtual Cluster* abstraction to model virtual networks with strict performance guarantees [5], [9], [26]. A virtual cluster offers the tenant the illusion for all her *Compute Units (CUs)* to be attached to a single non-oversubscribed switch with a minimum bandwidth  $b$  guaranteed. If excess bandwidth is available, it can be used in addition to the reserved bandwidth, e.g., leveraging recently proposed extensions to TCP such as Seawall [32].

A virtual cluster  $VC(n, b)$  has two parameters:  $n$ , the number of (identical) CUs in the cluster, and  $b$ , the bandwidth reservation from each CU to the virtual switch. Virtual clusters belonging to different tenants need to be embedded on a given *substrate*: a physical network connecting a set of servers. In this paper, we focus on multi-rooted tree (or *fat-tree*) like physical network topologies [1], [16] as they are the predominant topology in today's datacenters. These topologies are hierarchical and are recursively made of subtrees at each level. A fat-tree consists of a set of *pods* which are interconnected by *core routers*. Pods comprise a set of *racks* which are interconnected by the *aggregation switch*, and racks comprise multiple *servers* (or *hosts*) which are interconnected by the *Top-of-Rack (ToR) switch*. Each server can host a fixed number of CUs. As done in previous work, e.g., [5], [9], given the amount of multiplexing and assuming the availability of a multi-path routing protocol such as ECMP, we can approximate these links as a single aggregate link for bandwidth reservations.

To save costs, some datacenter operators introduce some degree of over-subscription, typically at the higher levels of the hierarchy. We model these configurations with two parameters  $\gamma_1, \gamma_2 \geq 1$  (called the *over-subscription factors* in [5]):  $\gamma_1$  denotes the factor of reduced capacity on the aggregation network (between ToR and aggregation switches) and  $\gamma_2$  the factor of reduced capacity between the aggregation switches and the core switch.

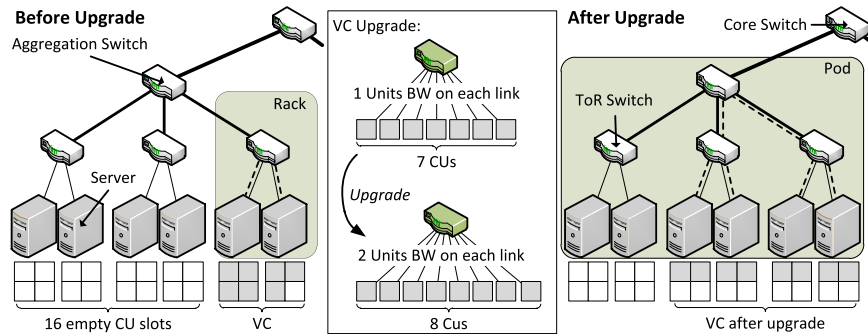


Fig. 3. Upgrade of a virtual cluster VC: *Left* - the initial state: VC(7,1) is embedded on the right-most rack of a pod of the fat-tree. The dashed lines indicate the current bandwidth reservations. *Middle* - the upgrade request: VC(7,1) needs to be upgraded to VC(8,2). *Right* - after the upgrade: Three CUs were migrated in order to find a new feasible embedding of VC which does not violate the capacity on the servers' uplinks.

The embedding of a virtual cluster describes its resource allocation in the substrate: an embedding maps each CU of the virtual cluster to a physical server in the substrate network; multiple CUs may be hosted on the same server. In addition, the embedding specifies the amount of bandwidth on each link reserved for the tenant. Intuitively, a “valid” embedding is one that does not oversubscribe server or network resources. A “good” embedding additionally chooses servers that are close in the physical network, thus minimizing unnecessary resource reservations on the physical links.

### B. The Challenge

The goal of this paper is to support virtual clusters whose guarantees can be adjusted over time, in an online fashion. Specifically, we want to be able to (1) upgrade a virtual cluster  $VC(n,b)$  consisting of  $n$  CUs and with a bandwidth guarantee  $b$ , both in *size* (i.e., number of CUs) as well as in the minimum *bandwidth*, that is, to a virtual cluster with  $x \geq 0$  more nodes and a factor  $\delta \geq 1$  more bandwidth, i.e., to  $VC(n+x, b \cdot \delta)$ ; (2) downgrade a virtual cluster in both size and bandwidth; (3) or a combination of both (e.g., upgrade size and downgrade bandwidth).

How to support such reconfigurations is also an algorithmic problem. Ideally, new feasible embeddings should be efficiently computable, i.e., at low runtime; moreover, we would like to avoid or at least minimize migrations in order to satisfy a reconfiguration request; finally, the resulting embeddings should have small network footprints, in the sense that no unnecessary bandwidth is reserved (on substrate links) to implement the virtual cluster guarantees.

### C. Example

To illustrate both the model and the challenge, let us consider an example. Figure 3 (*left*) shows a part of a fat-tree, i.e., a single pod consisting of three racks with two servers each; each server has 4 CU slots. We assume that the uplinks of the servers have a capacity of 4 units and the fat-tree provides full bisection bandwidth ( $\gamma_1 = \gamma_2 = 1$ ), resulting in a capacity of 8 units on the ToR switches' uplinks and a capacity of 24 units on the links between the aggregation switches and the core switch. On the right most rack, currently a virtual cluster

VC is embedded; the dashed line indicates the path along with bandwidth is reserved to connect the CUs. At some point, VC is upgraded, from VC(7,1) to VC(8,2), see Figure 3 (*middle*).

How can this request be satisfied? Theoretically, the right server in the rack still has a free CU slot which could be used to accommodate the additional CU; however, doubling the bandwidth reservations for each the CUs will violate the bandwidth capacities on the uplinks of the servers. Hence it becomes necessary to distribute the CUs in the substrate, in order to reduce the bandwidth utilization of the uplinks of the two servers. Thus, in this scenario, some CUs need to be migrated to satisfy the request. Figure 3 (*right*) shows a solution: the resulting embedding is valid.

## IV. THE SYSTEM

In this section, we first formalize the goals of the developed system, and then introduce the main concepts underlying Kraken and describe its key components.

### A. Objectives

Kraken is designed to accept and implement any embedding and upgrade request whenever there are sufficient resources available in the substrate. Downgrade requests, instead, can always be satisfied.

Besides satisfying upgrade requests whenever this is possible, Kraken is designed (1) to optimize the embedding cost of the virtual cluster, i.e., the amount of bandwidth which needs to be reserved in the physical network to host the virtual cluster; and (2) to reconfigure existing embeddings locally, i.e., to minimize the migration cost. To avoid affecting the performance of other tenants, we do not allow the migration of CUs belonging to other tenants, although in some cases this might result in lower embedding costs. The standard metric to evaluate the embedding cost (see also [5], [9]), is to measure the *embedding footprint*  $F(VC)$  of a virtual cluster VC:  $F(VC)$  is given by the overall network resources consumed by the VC, i.e., the sum of bandwidth reservations over all substrate links. (Note that the number of used CU slots is independent of the embedding.)

In order to measure the *reconfiguration costs*, we count the number of CUs which need to be embedded to a different location during an upgrade.

Notice that there is a trade-off between the two metrics: sometimes, at the price of higher reconfiguration costs, smaller footprints can be realized. In the following, we design our algorithms according to the following priorities (cf Section IV-F for a discussion of alternative objectives supported by Kraken): (1) the top priority is to satisfy a reconfiguration request; (2) the second priority is to minimize reconfiguration costs; and (3) the third priority, is to minimize the embedding footprint, i.e., among all solutions of the same reconfiguration costs, we compute the most resource efficient embedding.

Kraken provides the following worst-case guarantees.

- 1) *Request Satisfiability*: As long as a feasible solution exists all upgrade and downgrade requests are satisfied.
- 2) *Minimal Reconfiguration*: The reconfiguration cost is always minimized. In particular, if a solution without migrations exists, it is used. CUs of other tenants are never migrated.
- 3) *Optimal Allocation*: Among all possible solutions with minimal reconfiguration costs, Kraken computes the one with the minimal embedding footprint.
- 4) *Complexity*: The time complexity of re-configuring (or embedding) a virtual cluster is linear in the substrate size, in the worst-case.

### B. Algorithmic Concepts

At the heart of Kraken lie two main concepts: (1) The *center-of-gravity* (or simply: *center*) of a virtual cluster and (2) the *slotCount values*. The center-of-gravity concept (introduced in [31]) allows us to decouple the embedding of the individual Compute Units (CUs), in the sense that, given the location of the center-of-gravity, the CUs can be mapped “greedily”, one after the other, avoiding the combinatorial complexity and rendering the problem polynomial time solvable. The *slotCount(v)* values provide an aggregate information about the number of available CU slots in the sub-tree of the fat-tree below a given node  $v$ ; they constitute the main data structure used by Kraken. While previous virtual cluster embedding algorithms used a similar concept [5], [9], [14], only the combination with the center-of-gravity concept allows a modification which enables the low runtime of the dynamic algorithm (roughly linear in the substrate size).

**Center-of-Gravity.** The virtual cluster abstraction offers tenants a network where each CU is connected to a virtual switch at bandwidth  $b$  [5]. While this virtual switch is only a logical concept, its position in the substrate matters, as resources need to be reserved from it to each CU.<sup>1</sup> The center-of-gravity may also be located on a server, not only on a switch (e.g., if many CUs of the virtual cluster are collocated on the same server). Given a mapping of the CUs of a given virtual cluster VC, we will refer to the optimal position of the virtual switch (with respect to embedding footprint) as the *center-of-gravity* CoG of VC.

Given any node  $v$  in the fat-tree (either a server or a switch), we can partition the nodes of VC into two sets with respect

to  $v$ : the set of CUs *at or below* the node  $v$  in the fat-tree, and the remaining CUs *above* (or “outside”)  $v$ . Sometimes, we use the same terminology to refer to the location of substrate components relative to each other.

When applying the COG concept to the fat-tree topology, we have two important properties, which Kraken leverages: (1) no more than half of the nodes, can be *above* CoG and (2) no more than half of the nodes are *below* one of the children of CoG. The correctness of this property can be shown easily by contradiction: If more than half of the CUs are behind one link, moving the CoG in this direction will decrease the bandwidth costs for more than half of the CUs by 1 and increase the costs for the other CUs by 1, resulting in a smaller footprint.

Moreover, when computing the embedding footprint of a virtual cluster VC, it is often helpful to count the number of CUs which are embedded *below* CoG(VC); we will refer to this number as  $\beta$ . The remaining CUs of VC which are embedded *above* CoG(VC), fall into three classes: the  $\alpha^{(p)}$  “far-away” CUs located in a different pod, the  $\alpha^{(r)}$  CUs in the same pod but in a different rack, and the  $\alpha^{(s)}$  CUs in the same rack but on a different server. This classification results in simple formulas for the embedding footprint of a virtual cluster. For instance, if CoG(VC) is embedded to a top-of-rack switch, the embedding footprint is given by  $F(\text{VC}) = \beta + 3 \cdot \alpha^{(r)} + 5 \cdot \alpha^{(p)}$  as the distance to servers in the same rack ( $\beta$ ) is 1 and the distance to all servers in the same pod but in different racks ( $\alpha^{(r)}$ ) is 3 while the distance to servers in other pods ( $\alpha^{(p)}$ ) is 5.

**slotCount-Values.** The second core concept of Kraken is the *slotCount(v)-value*: intuitively, the *slotCount(v)-value* indicates how many additional CUs can be placed below a certain substrate node  $v$  (a server or switch), such that the currently available server and link resources are all satisfied.

The number of CUs which can be placed below a certain substrate node  $v$  depends on two factors: the available bandwidth and the available CU slots. For Kraken it is sufficient to compute the bandwidth criteria for cases where CoG is *above*  $v$ . This eases the computation of these values significantly, since the resulting interval of possible amounts of CUs becomes continuous. In order to keep the runtime of the *slotCount* computation low, we leverage the optimal sub-problem property in our dynamic program: We start by computing the *slotCount*-values on the host level. For each server  $s$  we compute  $\text{slotCount}(s) = \min(\text{spareCUs}(s), \lfloor \text{spareBW}(s) / b \rfloor)$  where  $\text{spareCUs}(s)$  denotes the available CU slots of a server  $s$  and  $\text{spareBW}(s)$  denotes the available bandwidth on the uplink. The *slotCount* of a rack  $r$  is then defined as:  $\text{slotCount}(r) = \min(\sum_{s \in r} \text{slotCount}(s), \lfloor \text{spareBW}(s) / b \rfloor)$ . The *slotCount(p)*-values for pods can subsequently be computed from the racks’ *slotCount*-values.

**Overview.** Based on these concepts, in order to embed or reconfigure a virtual cluster VC, Kraken simply cycles through all possible center-of-gravity locations in the substrate network (servers and switches): for each possible CoG location  $v$ , Kraken determines the minimal number of migrations needed, in order to shift the center to  $v$ . This is a fast operation since it does not scale with the size of the substrate, but with the size

<sup>1</sup>Note that there could be multiple positions with the same embedding cost, and that in a fat-tree, a distributed switch mapping does not reduce costs.

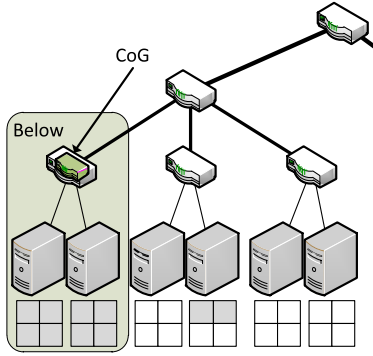


Fig. 4. Example embedding of a  $VC(10,1)$ . The CoG is on the left top-of-rack switch. The 8 CUs located in the left most slot are *below* the CoG and entail a unit bandwidth cost each. The two remaining CUs are *above* ( $\alpha^r$ ) the CoG and inflict three units of bandwidth cost each. Moving the virtual switch to one of the servers in the left rack would reduce the bandwidth costs for the 4 CUs on that server by 1, but increase the costs for the other 6 CUs by 1. Moving it to the pod would decrease the costs for the two CUs in the middle rack, but increase the costs for the other 8 CUs.

of the VC. If CoG can be implemented on  $v$  with minimal migration costs, the *slotCount* values are used to calculate the best possible embedding footprint of a mapping with the center at  $v$ . As we will show, this also does not require scanning the entire substrate, and is fast. For an illustration, we refer to Figs. 4 and 5.

### C. Upgrade Algorithm

Algorithm 1 shows the pseudo-code of Kraken’s algorithm to implement an upgrade operation `upgrade`, from  $VC(n, b)$  to  $VC(n + x, \delta \cdot b)$  with  $x \geq 0$  more nodes and a factor  $\delta \geq 1$  more bandwidth. We use  $\mu$  to denote the embeddings.

---

#### Algorithm 1 Algorithm Upgrade( $VC, x, \delta$ )

---

*Output:* success or failure

```

1: for all nodes  $v$  in the fat-tree: compute  $slotCount(v)$  values
2:  $m^* \leftarrow \infty$ ;  $F^* \leftarrow \infty$ ;  $cog^* \leftarrow \perp$ ;
3: for all  $v$  in substrate do
4:    $M \leftarrow \text{minMig}(v)$ 
5:   if  $|M| \leq m^*$  then
6:      $F \leftarrow \text{footprint}(v, |M|)$ 
7:     if  $F < \infty \wedge (|M| < m^* \vee F < F^*)$  then
8:        $cog^* \leftarrow v$ 
9:        $m^* \leftarrow |M|$ 
10:       $F^* \leftarrow F$ 
11:    end if
12:  end if
13: end for
14: if  $m^* = \infty$  then
15:   return failure
16: end if
17:  $\mu \leftarrow \text{computeEmbedding}(VC, cog^*)$ 
18: return success
    
```

---

Kraken first pre-computes the *slotCount*-values for the entire substrate network, i.e., for each substrate node  $v$  (a server or switch). Subsequently, Kraken computes the new center-of-gravity CoG for VC which minimizes the reconfiguration

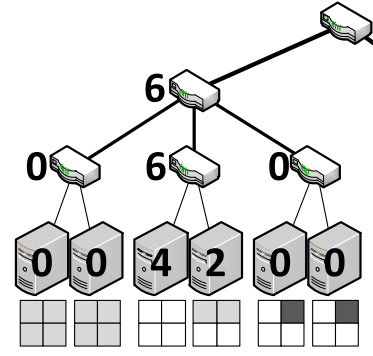


Fig. 5. Example for *slotCount* computations. Already embedded are a  $VC_1(10,1)$  (light gray) and a  $VC_2(2,4)$  (dark gray). We assume that the substrate has a maximum bandwidth of 4 on the links connecting the servers to the top-of-rack switches and that the currently processed request has the form  $VC_3(n,1)$ . Obviously, the left rack is full and hence all slot count values are 0. The center rack has one completely free server (*slotCount* = 4) and one half free server (*slotCount* = 2) resulting in a *slotCount* = 6 on the top-of-rack switch. The two servers in the right rack are partially occupied by  $VC_2$ , but could potentially host 3 CUs of  $VC_3$  (e.g., when  $n = 3$ ). However, *slotCount* is based on the assumption, that CoG is *above* the node and the two CUs of  $VC_2$  occupy the entire bandwidth. Hence, the *slotCount* values are 0 for this rack, too.

costs in terms of the number of to be released, i.e., migrated CUs  $M$  (function `minMigs`) and embedding footprint  $F$  (function `footprint`), by iterating over all nodes in the substrate. Subsequently, the best found solution is embedded (function `computeEmbedding`).

1) *Minimal Migrations*: To compute the minimal number of migrations, function `minMig` proceeds as follows, see Algorithm 2: For each node  $v$  in the substrate (i.e., all servers and switches), it computes a list of CUs which have to be “released” (i.e., put in a pool of CUs which will be embedded somewhere else by the algorithm), to be able to realize the new center-of-gravity at node  $v$ .

---

#### Algorithm 2 Minmig(Substrate Node $v$ )

---

*Output:* set of CUs

```

1:  $M \leftarrow \emptyset$ 
2:  $L \leftarrow \text{computeConflictLinks}(v)$ 
3: sort  $L$  with decreasing distance from  $v$ 
4: for all links  $\ell \in L$  do
5:   while  $\ell$  oversubscribed do
6:     let  $c$  be an arbitrary CU below  $\ell$ 
7:      $M \leftarrow M \cup \{c\}$ 
8:   end while
9: end for
10:  $M \leftarrow M \cup \text{extraCUs}(v)$ 
11: return  $M$ 
    
```

---

`ComputeConflictLinks` computes the set of links  $L$  whose capacity would be oversubscribed if the center-of-gravity  $cog$  was on  $v$  and the bandwidth was increased to  $b \cdot \delta$  under the current embedding  $\mu$  of the existing CUs. Subsequently, we iteratively release CUs until a critical link  $\ell \in L$  is no longer oversubscribed. This yields the first part of the set  $M$  of CUs which need to be migrated. The conflict resolution is ordered by distance to the center-of-gravity.

While releasing the CUs so far in  $M$  ensures that no link is oversubscribed, additional CUs may have to be moved to guarantee that the center-of-gravity is realized at the desired physical node: thus, `extraCUs` adds more CUs to the set  $M$ , such that the sum of the CUs which are currently hosted below  $v$  and the cardinality of  $M$  reach  $n/2$ . To make  $v$  the center-of-gravity of the virtual cluster, it is necessary and sufficient that at least  $n/2$  CUs are below  $v$ .

2) *Minimal Footprint*: After determining the number of CUs that have to be migrated, we compute the embedding footprint. Interestingly, Kraken can compute the embedding cost of a desired center-of-gravity *without* determining an explicit embedding of the new virtual cluster, by utilizing the `slotCount`-values.

The function `footprint` is described in Algorithm 3. It takes a desired center-of-gravity  $v$  and a target number  $m$  of CUs which are to be migrated. Let us first observe that the footprint of a virtual cluster can be computed via the following case distinction: (1) If  $v$  is a core switch, all CUs are located below  $v$ , and hence the distance between  $v$  and the CUs is three. Thus,  $F(\text{VC}) = 3 \cdot \beta$ , where  $\beta$  counts the number of CUs which are embedded *below*  $\text{CoG}(\text{VC})$ . (2) If  $v$  is an aggregation switch of a pod, the CUs of VC are either located on servers in the same pod, or on servers in different pods. Clearly, all servers in the same pod are at distance two from  $v$ , and the servers in other pods are at distance four from  $v$ . We have  $F(\text{VC}) = 2 \cdot \beta + 4 \cdot \alpha^{(p)}$ , where  $\alpha^{(p)}$  is the number of CUs of VC which are embedded *above*  $\text{CoG}(\text{VC})$ , in a different pod. (3) In case  $v$  is embedded to a ToR switch, the embedding footprint is given by  $F(\text{VC}) = \beta + 3 \cdot \alpha^{(r)} + 5 \cdot \alpha^{(p)}$ , where  $\alpha^{(r)}$  is the number of CUs of VC which are embedded *above*  $\text{CoG}(\text{VC})$ , in a different rack. (4) The embedding footprint for a  $v$  on servers is given by  $F(\text{VC}) = 2 \cdot \alpha^{(s)} + 4 \cdot \alpha^{(r)} + 6 \cdot \alpha^{(p)}$ , where  $\alpha^{(s)}$  is the number of CUs of VC which are embedded *above*  $\text{CoG}(\text{VC})$ , on a different server. In this case, CUs which are embedded *below* the COG are omitted, as they have no bandwidth costs.

---

**Algorithm 3** `Footprint`(Substrate Node  $v$ , Number of CUs to Migrate  $m$ )

---

*Output*: cost value

```

1: done  $\leftarrow$  0
2: for all children  $v'$  of  $v$  in the fat-tree do
3:   done  $\leftarrow$  done + slotCount( $v'$ )
4: end for
5: return  $\text{ST}(v) + \text{height}(v) \cdot n + \text{costsAbove}(v, m - \text{done})$ 

```

---

The function `footprint` first computes the number of CUs which can be placed on each of the sub-trees represented by the direct children of  $v$ . Since the center-of-gravity  $v$  is above its children by definition, the `slotCount`( $v$ )-values of the children are accurate. Then, the embedding cost is computed recursively by the formula  $\text{ST}(v) + \text{height}(v) \cdot n + \text{costsAbove}(v, z - \text{done})$ . The first cost term  $\text{ST}(v)$  accounts for the static costs, i.e., the costs from CUs which are not scheduled for migration according to the minimal migrations. The second cost term  $\text{height}(v) \cdot n$  depends on the depth of

the center-of-gravity in the tree. The third term computes the additional costs from the CUs above  $v$ , if any, see the function `costsAbove` (Algorithm 4): we leverage the fact that the costs for placing CUs further away from a candidate center  $v$  increases by two for every layer in the fat-tree, regardless of the layer where  $v$  is located. Accordingly, given  $z$  flexible CUs, we add  $2z$  to the costs and execute the function again with the parent node of  $v$  as the new  $v$  and  $z - \sum_{v' \in V'} \text{slotCount}(v')$  as the new  $z$ , where  $V'$  is the set of siblings of  $v$  (i.e., children of the parent node of  $v$  excluding  $v$ ). If  $v$  is the core switch, or the spare capacity on the uplink of  $v$  is less than  $z \cdot \delta \cdot b$ ,  $v$  cannot be the center-of-gravity, and the upgrade request fails for this specific location of the COG. If this is the case for all nodes  $v$  in the substrate, the upgrade request has to be rejected.

---

**Algorithm 4** `CostsAbove`(Substrate Node  $v$ , Number of Flexible CUs  $z$ )

---

*Output*: cost value

```

1: if  $z = 0$  then
2:   return 0
3: end if
4: if ( $v$  is a core switch or the uplink from  $v$  does not have
    $z \cdot \delta \cdot b$  spare bandwidth) then
5:   return  $\infty$ 
6: end if
7: done  $\leftarrow$  0
8: for all for all siblings  $v'$  of  $v$  do
9:   done  $\leftarrow$  done + slotCount( $v'$ )
10: end for
11: return  $2 \cdot z + \text{costsAbove}(\text{parent of } v, z - \text{done})$ 

```

---

#### D. Downgrade Algorithm

Downgrade operations in Kraken never require any migrations. However, the center-of-gravity may change. Thus, the downgrade algorithm of Kraken proceeds similar to the upgrade algorithm, but without functions `minMig` and without the need to compute the `slotCount`( $v$ ) values. The main difference regards how the values are actually used to compute the costs. While the original algorithm depends on `slotCount`-values and the current distribution, we set the current distribution to 0 and all `slotCount`-values to the distribution prior to the upgrade.

#### E. Formal Guarantees

Since the calculated cost and `slotCount` values are *exact*, we have derived the following result.

*Theorem 1: Kraken guarantees:*

- 1) Request Satisfiability: *As long as a feasible solution exists all upgrade and downgrade requests are satisfied.*
- 2) Minimal Reconfiguration: *The reconfiguration costs is always minimized. In particular, if a solution without migrations exists, it is used.*
- 3) Optimal Allocation: *Among all possible solutions with minimal reconfiguration costs, Kraken computes the one with the minimal embedding footprint.*

- 4) Complexity: *The time complexity of re-configuring (or embedding) a virtual cluster is bounded by  $O(N \cdot n \cdot \Delta)$  in the worst-case, where  $N$  is the size of the substrate (number of servers),  $n$  is the virtual cluster size, and  $\Delta = S + R + P$  is the number of servers in a single rack  $S$  (i.e., the degree of a ToR switch), plus the number of racks in a single pod  $R$  (i.e., the degree of an access switch), plus the number of pods  $P$  (i.e., the degree of a core switch).*

*Proof:* The optimality proof unfolds in two central lemmas: Lemma 1 is the key to the *Minimal Reconfiguration* property, and Lemma 2 is the key to the *Request Satisfaction* property. Finally, in Lemma 3, we prove the runtime complexity.

*Lemma 1:* *The function `minMigs` computes the minimal number of migrations required to shift the center-of-gravity to a given substrate node  $v$ .*

*Proof:* We first note that a node  $v$  can only represent a center-of-gravity of a virtual cluster VC under a new embedding  $\mu'$  if (1) in  $\mu'$  at least  $n/2$  CUs are embedded below  $v$ , and if (2)  $\mu'$  describes a feasible embedding, i.e., no nodes or links are oversubscribed (or *critical*). If Condition (1) is violated, the allocation cost under an alternative center (*above*  $v$ ) is strictly lower.

Let  $L$  be set the of conflict links of embedding  $\mu$  (function `computeConflictLinks`). For each conflict link  $\ell \in L$ , we can define the non-empty set  $S_\ell$  of CUs whose bandwidth reservation involves  $\ell$ , and whose removal will reduce the load of link  $\ell$ . Due to the hierarchical structure of the substrate, if two links  $\ell$  and  $\ell'$  are on the same path from a host to center-of-gravity  $v$ , and  $\ell$  is lower than  $\ell'$  in the tree, it holds that  $S_\ell \subseteq S_{\ell'}$ . Thus, removing CUs in a descending order of distance to  $v$  (see function `minMigs`), will minimize the number of necessary changes. Finally, `minMigs` fulfills Condition (2) by adding only the necessary number of CUs below  $v$ .  $\square$

Since Kraken iterates over all possible center-of-gravity locations  $v$ , Lemma 1 directly implies the *Minimal Reconfiguration* property.

The next lemma shows that Kraken will always find a feasible realization of a center-of-gravity  $v$  if it exists.

*Lemma 2:* *Function `footprint(v, |minMig(v)|)` only returns  $\infty$  if it is unfeasible to make  $v$  the center-of-gravity.*

*Proof:* From the proof of Lemma 1 we know that no solutions exist with  $k < |M|$  migrations, where  $M$  is as defined in `minMig`. Function `footprint` first computes the embedding costs at and below  $v$ , after which  $x = |M| - \sum_{v'} \text{slotCount}(v')$  conflicts below  $v$  are left. While each additionally released CU below  $v$  frees up one slot for some child  $v'$ , it at the same time increases the set of conflicts from 1 to  $|M|$ . Thus, the resulting number of conflicted CUs to place *above*  $v$  remains  $z$ . A cost `footprint(v, |minMig(v)|) = \infty` implies that  $z > 0$ , since for  $z = 0$ , `costsAbove(v, 0) = 0` and hence `footprint(v, |minMig(v)|) < \infty`.

On each tier in the substrate, `costAbove(v', z')` can return  $\infty$  if the uplink of  $v'$  does not have enough spare bandwidth ( $b \cdot \delta \cdot z'$ ). Since `minMig` generated a conflict-free partial embedding, we get  $z' > 0$  if there is a conflict and  $b$  and  $\delta$  are

positive. The relation  $z' > 0$  implies that no resources are left on the previous tiers ( $z'' > \sum \text{slotCount}(v'')$ ). Hence, each released CU will increase the spare capacity on the uplink by  $b \cdot \delta$ , but also increase  $z'$  by one, inhibiting a feasible solution. Once function `costAbove` reaches the core switch,  $\infty$  is returned. In this case, the sum of all `slotCount(v)`-values for the pods is less than  $|M|$ . Releasing a CU in any pod will increase the `slotCount(v)`-values for that pod by one and at the same time increase  $|M|$  by one. Hence this conflict cannot be resolved by migrating additional CUs.  $\square$

Again, since Kraken iterates over all locations  $v$ , the *Request Satisfaction* holds. The *Optimal Allocation* property can be shown along the same lines. The time complexity of Kraken is as follows.

*Lemma 3:* *The time complexity to satisfy a request is bounded by  $O(N \cdot n \cdot \Delta)$  in the worst-case, where  $N$  is the size of the substrate (number of servers),  $n$  is the virtual cluster size, and  $\Delta = S + R + P$  is the number of servers in a single rack  $S$  (i.e., the degree of a ToR switch), plus the number of racks in a single pod  $R$  (i.e., the degree of an access switch), plus the number of pods  $P$  (i.e., the degree of a core switch).*

*Proof:* The computation of the `slotCount(v)`-values requires  $O(N \cdot \Delta)$  time as the dynamic program runs in a bottom-up manner. Subsequently, Kraken iterates over all possible center-of-gravity locations in the fat-tree (time  $O(N)$ ): for each candidate COG, conflicts are computed along the links from the assigned server of each CU to the potential COG (time  $O(n)$ ), together with the costs for the resulting embedding (time  $O(\Delta)$ ). The overall runtime for finding the optimal COG is hence  $O(N \cdot n \cdot \Delta)$ . The actual embedding can utilize the previously computed list of conflicts and distribution of migrated CUs across the substrate. To generate a feasible embedding from here it is necessary to traverse through the sub-trees which should host the migrated CUs afterwards. During this traversal, we visit at most  $N$  nodes and check their `slotCount(v)` values, again at a runtime of  $O(N)$ . Hence, the overall runtime amounts to  $O(N \cdot n \cdot \Delta)$ .  $\square$

Note that Kraken can also be used to embed virtual clusters from scratch, and ensuring a minimal footprint. Thus, together with property 4), Kraken also outperforms state-of-the-art virtual cluster embedding algorithms which do not support any reconfigurations, e.g., [14], [31], at least in the worst-case: Our simulations show that Oktopus [5] and Proteus [9] find fairly good embeddings with small footprints as well. However, in the worst case, their performance can be arbitrarily bad compared to Kraken. In the case of Oktopus, a small virtual cluster which could be hosted by a single server (Kraken footprint: 0) may be embedded across multiple servers (footprint  $> 1$ ). In the case of Proteus, the ratio of the optimal footprint computed by Kraken and the footprint by Proteus can be as high as  $n/3$ : such an example can be constructed by exploiting the fact that Proteus will only consider cross-pod embeddings if a request cannot be fit in a single pod. Thus, in case  $n - 1$  slots are available on a single server in one pod and  $n$  times one slot is available on servers in another pod, the Proteus footprint is  $2n$  while the footprint of Kraken is 6.



### F. Alternative Migration Cost Models

For ease of exposition, we presented Kraken for a simple model where the objective of minimizing the number of migrations is prioritized over optimizing the embedding footprint. However, our algorithms can be extended to other migration cost models and trade-offs between migration and footprint costs, without sacrificing optimality. For instance, intra-pod migration costs could be modeled to be cheaper than inter-pod migrations, and migration costs could also depend on the available bandwidth along the migration path.

## V. EVALUATION

We conduct extensive simulations to study the feasibility of online reservation upgrades at runtime. By default, we will assume the same settings and parameters as used in previous work [5]. However, given our more dynamic environment, we also introduce a model for elastic reconfiguration requests, and conduct a sensitivity analysis, studying the impact of different factors (such as magnitude of reconfiguration and system load) by using parameter sweeps.

### A. Metrics

We consider the following two metrics:

**Acceptance Ratio.** Ideally, a system such as Kraken should be able to accept and satisfy as many requests as possible. For each request (either arrival of a new virtual cluster or a reconfiguration request), we distinguish whether or not the request was satisfied and, if satisfied, whether it was satisfied (1) *with* or (2) *without migrations*. Note that Kraken does not use “strategic access control” (e.g., to favor “small” requests to improve that acceptance ratio); in fact, Kraken never rejects a request if it can be satisfied.

**Reconfiguration Costs.** While our simulation does not capture many parameters that determine the actual cost of a migration, we count the number of migrations; this is a natural metric given the uniform size of CUs of the virtual cluster. In particular, we will report on the *fraction* of migrated CUs relative to the virtual cluster size, which provides more insights than an absolute number.

### B. Methodology & Runtime

**Substrate.** We model the datacenter as a three-level fat-tree. Overall, we have 16,000 servers distributed over  $P = 10$  pods of  $R = 40$  racks each; a rack contains  $S = 40$  servers. By varying the connectivity and the bandwidth of the links between the switches, we change the over-subscription of the physical network. By default, we will assume that the access network is oversubscribed by a factor  $\gamma_1 = 4$ , while the core is not oversubscribed ( $\gamma_2 = 1$ ). The available bandwidth is  $B = 10$  Gbps.

**Demand.** New virtual cluster requests arrive according to a Poisson process with  $\lambda = 0.36$ . The lifetime of each virtual cluster is chosen according to an exponential distribution with average 3,600 s (one hour). By default, the size of a virtual cluster and the bandwidth are chosen from an exponential distribution with mean 49 and 2.5 Gbps respectively.

The parameters are normalized to induce a system load of 0.75 on average. The size of the virtual cluster in numbers of CUs is chosen randomly from an exponential distribution, with an average of 49 CUs per cluster.

**Elastic Model.** To add dynamicity to the virtual cluster demands, we use six additional Poisson processes<sup>2</sup> which continuously pick virtual clusters for upgrading and/or downgrading in a *multiplicative manner*. More precisely, the embedded clusters are continuously reconfigured by these six independent processes which randomly choose one of the existing clusters and perform a multiplicative update, i.e., either (1)+(2) upgrade or downgrade the bandwidth by a factor  $f_b$  ( $f_b$  corresponds to  $\delta$  in our formal sections), (3)+(4) increase or decrease the cluster size by a factor  $f_n$  ( $f_n$  is the multiplicative version of the additive  $x$  in our formal sections), (5)+(6) *jointly* upgrade or downgrade the bandwidth *and* the cluster size by a factor  $f$ . By default, we assume that  $f = f_b = f_n = 1.5$ . With regards to reporting the results, we focus on the upgrades as these are the ones which trigger migrations.

To ensure the statistical significance, we run our simulations for 80,000 rounds, which is roughly eighty times the duration (i.e., lifetime) of a virtual cluster. To avoid artifacts related to the initial empty substrate, we omit the first 10k requests.

**Runtime.** In this scenario, Kraken requires 86 ms on average to satisfy any given request (the 99<sup>th</sup> percentile is 344 ms), when run on an Intel i3-2310M CPU @ 2.10GHz.

### C. Baseline Comparison

Kraken features two main mechanisms for the efficient upgrade of a virtual cluster: (1) Kraken allows to upgrade an existing embedding by increasing the bandwidth between CUs *at their current locations*, as well as by the extending the cluster by the local addition of new CUs; (2) if a local extension is not sufficient to satisfy a request, Kraken also supports the re-embedding, i.e., migration of existing CUs.

In order to understand the contribution of each of these two features, we break down the analysis of Kraken into two steps: We first study a variant of Kraken, called KrakenNP, which does not perform fine-grained migrations. (*NM* stands *No (local) Migrations*.) That is, KrakenNP is equivalent to Kraken, but if a request cannot be satisfied with the given CU embedding, it resorts to embedding the virtual cluster with the new specification from scratch. Subsequently, we study the full-fledged Kraken system which can migrate CUs arbitrarily in order to satisfy requests (subject to the usual constraint that the number of migrations should be kept minimal). For a simple baseline comparison, we also re-implemented Oktopus [5]; we extended Oktopus so that requests can be satisfied by re-embedding.

To give a basic understanding of the number of migrations required to support elastic virtual clusters, Figure 6 plots the empirical cumulative distribution function (ECDF) of the migration cost for the three algorithms KrakenNP, Kraken and Oktopus, and the three operations: add CUs,

<sup>2</sup>While Poisson distributions are commonly used to describe arrival patterns, we still lack good empirical models for the kinds of workloads considered in this paper.

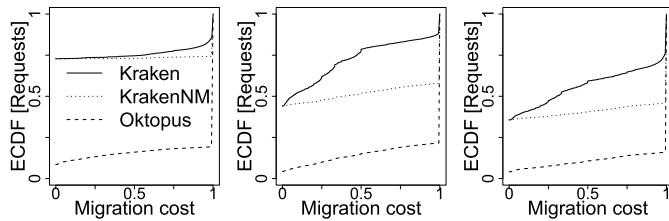


Fig. 6. Reconfiguration costs: KrakenNP vs. Kraken vs. Baseline (augmented Oktopus)—(left:) cluster size upgrade, (middle:) bandwidth upgrade, (right:) joint upgrade. The legend on the left is valid for all three plots.

upgrade bandwidth, and joint upgrade of CUs and bandwidth. Note that when a new embedding is performed to satisfy an upgrade request, the mechanism will guide the embedding process to a similar configuration. This means that when possible, the CUs will be assigned to the same old location, which, hence, will not be counted toward the migration cost. This explains why in some cases the migration cost of Oktopus and KrakenNP can also have values different from zero (no migrations) and one (all CUs are migrated).

We first discuss a scenario where only the bandwidth is upgraded. In Figure 6 (middle), we can observe that already KrakenNP is far superior to Oktopus as it can satisfy 45% of the upgrade requests without migrations at all, while Oktopus has to migrate all CUs of a VC for 80% of the upgrade requests. In general, we find that Oktopus will likely find similar embeddings (with few migrations) if the upgrade request happens temporally close to the embedding time. However, later it becomes likely that virtual clusters will be embedded on a different sub-tree (or pod), resulting in many migrations. The performance of Kraken is very similar to the one of KrakenNP. However, the missing support of partial and coordinated migrations leads to  $\approx 50\%$  cases where KrakenNP has to migrate all CUs, while Kraken can avoid migrating more than 50% of the CUs for nearly 80% of the requests.

The corresponding results for cluster size upgrades are shown in Figure 6 (left). While Oktopus can only embed about 10% of the upgrade requests without migrating any CUs, Kraken can upgrade 70% of the requests without migration. KrakenNP achieves a similar performance, and only for 10% of the requests, we can observe an improvement  $\geq 5\%$  with Kraken in terms of reconfiguration costs.

Figure 6 (right) studies joint upgrades (bandwidth and cluster size). Here, the overall performance of Oktopus remains the same, and the performance of Kraken and KrakenNP becomes a mixture of the previous cases. While both variants of Kraken need no migrations for 35% of the requests, KrakenNP has to migrate all CUs for 40% of the requests, while Kraken can satisfy about 70% of all requests without migrating all CUs.

D. Sensitivity Study

Next, we conducted a sensitivity study of Kraken, in which we performed parameter sweeps for the up- and downgrade ratios  $f_b$  and  $f_n$ , the mean number of CUs per request, the bandwidth requirements per CU, the substrate load, and the access network over-subscription ratio. We will first study

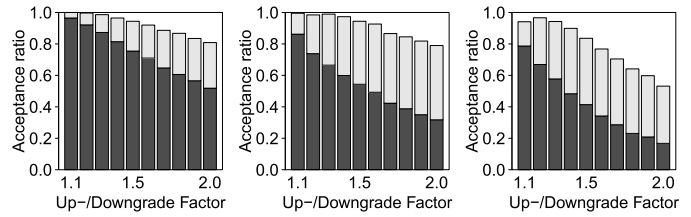


Fig. 7. Kraken acceptance ratios: without migration (dark gray), with migration (light gray)—(left:) cluster size upgrade, (middle:) bandwidth upgrade, (right:) joint upgrade.

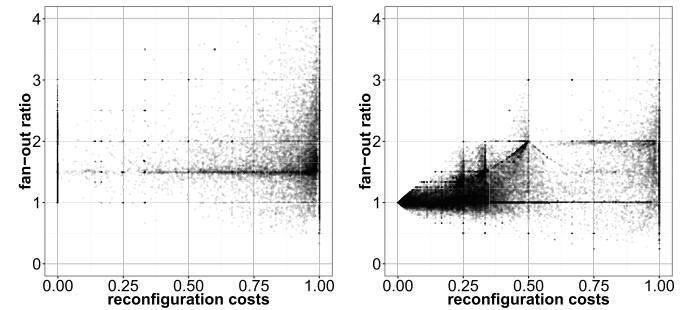


Fig. 8. Scatter plot of reconfiguration cost vs. fan-out ration—(left:) cluster size upgrade, (right:) bandwidth upgrade.

the effect of the upgrade ratios  $f_b = f_n$  in greater detail, and subsequently, we report on our general observations for the other parameters.

Figure 7 shows the acceptance ratio for virtual cluster upgrades as bar plots. The dark gray area corresponds to upgrade requests that do not require migration. The light gray component of the bar corresponds to those requests that can be satisfied by Kraken but require migration. We again have three subplots corresponding to the three operations: adding CUs, upgrading bandwidth, and joint upgrades of CUs and bandwidth. The impact of the upgrade factor  $f$  is significant, opening a spectrum from “accepting almost all requests without migrations” (for factors close to one) to “no migration for only 50% of the cluster size upgrade requests”. The impact of  $f$  on the bandwidth upgrades is even more articulated. As expected in the joint upgrade scenario, the two factors are amplified. Indeed, the problem is unfeasible for more than 40% of the requests if the upgrade factor is 2.

To explain why it is often possible to satisfy cluster size upgrade requests without migrations, Figure 8 plots (for upgrade factor 1.5) the reconfiguration cost (in terms of migrations) vs. the fan-out ratio (number of used servers), i.e., the fan-out after the upgrade operation divided by the fan-out before the upgrade. Let us start with the cluster size upgrades. Figure 8 (left) has two interesting areas: First, the one with no migrations which corresponds to zero reconfiguration costs. Here the fan-out ratio increases up to a factor of three as additional CUs can simply be added in the corresponding sub-tree. Second, those that require migrations require full re-configurations and therefore fall to the right hand side of the plot. There are not that many requests in between the two extremes, and the fan-out is often around 1.5 if there

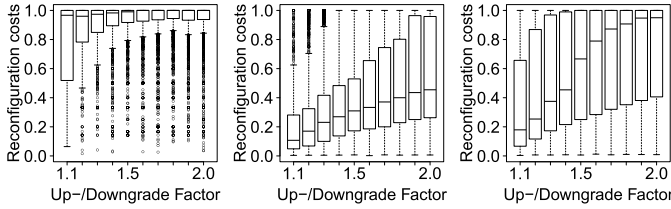


Fig. 9. Kraken reconfiguration costs for upgrades with migrations—(left:) cluster size upgrade, (middle:) bandwidth upgrade, (right:) joint upgrade.

are partial migrations, which is reasonable given the 50 % upgrade. The plot for bandwidth upgrades (Figure 8 (right)) is quite different. Since only the bandwidth is upgraded, the fan-out does not change if no CU is migrated. In this case, there are many upgrade requests with relatively small reconfiguration costs while the cluster size upgrades created more extremes. This can be explained by the collocation strategy used by Kraken: for instance, if three CUs are hosted on a single server prior to the upgrade, a bandwidth increase may require that one CU be moved to another server to alleviate the load on the uplink.

To better understand the difference between adding CUs and upgrading the bandwidth, Figure 9 zooms into the light-gray area and plots the distribution of the relative number of migrations, given that the upgrade required at least one migration. While in most cases it is sufficient to migrate less than half of the CUs for bandwidth upgrades, it is necessary to migrate more than 90% of the CUs, if any reconfigurations are necessary during a size upgrade. This can be explained by the different triggers of migrations for the two operations: In many situations, the CUs of a VC are collocated with each other. Adding CUs in this cases does not require reconfigurations, as long as there is sufficient spare bandwidth on the subtree, which currently hosts the VC. Contrary, even a small bandwidth upgrade can change the maximum number of CUs which can be collocated (e.g., a bandwidth upgrade from 2.4 Gbps to 2.6 Gbps changes maximum number of collocated CUs from 4 to 3), which will require a share of the CUs (in this case 25%) to be migrated. The only case in which adding CUs will actually trigger migrations, occurs when the sub-tree which currently hosts the VC is already highly filled, and the center has to be moved in order to meet the bandwidth guarantees. This can also happen during a bandwidth upgrade, but the first case occurs more often, and hence has a strong impact on the outcome shown in Figure 9. The joint upgrade case, shows the combined effects of the other two described upgrades.

We will now report on our observations for the other parameters: Varying any of the above parameters by 50% never caused the acceptance ratio to drop below 80%. Moreover, the acceptance ratio for CU as well as bandwidth upgrades are comparable to those of Figure 7. Joint upgrades are slightly more complex but the acceptance ratio is still above 80%. The largest difference we observed in the worst case acceptance ratio was 6%.

With regards to the reconfiguration costs we find that cluster size upgrades are typically more expensive. This is fully consistent with the observations above. It also points out that

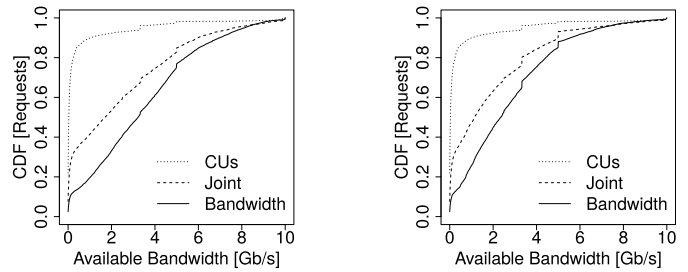


Fig. 10. CDF of the available bandwidth to migrate a compute unit for upgrades which require migrations. Left: avg. bandwidth; Right: min. bandwidth.

even local greedy search strategies for re-embedding CU size upgrades can be fairly successful.

In general, we see that most parameters only have a very small effect on the reconfiguration costs of bandwidth upgrades, and a small effect on the joint upgrade. On average across all evaluated parameters, bandwidth upgrades need approximately one third reconfigurations per CU, while joint upgrades typically require two third reconfigurations per CU. This indicates that these operations benefit from the rigorous optimizations of Kraken.

#### E. Bandwidth for Migrations

While compute units can be small and light-weight, it may sometimes be desirable to migrate more state or entire VMs. Therefore, we investigate the bandwidth available during CU migrations. Figure 10 shows that for bandwidth upgrades, on average, approximately 3 Gbps can be guaranteed along the migration path of each CU on average; the minimum is around 2 Gbps. For joint upgrades, the values are 2 Gbps on average and 1 Gbps for the CU with the lowest available bandwidth. These values are encouraging, indicating that even large migrations are feasible in reasonable time. However, we also see that on the occasion where cluster size upgrades trigger migrations, the bandwidth can become critical: only 10% of the requests can guarantee more than 1Gbps of bandwidth for the migrations. In such settings, one may have to resort to a separate management network for migration.

## VI. CASE STUDY WITH HADOOP

We now go back to the Hadoop case study introduced in Section II. Hadoop is an appealing application for Kraken, as the Hadoop framework natively supports tracking the progress of a task. In Hadoop-YARN, MapReduce tasks inform the Application Master about the progress of a task periodically. The Application Master uses this information, for instance, to speculate tasks that are straggling. A Kraken scheduler can easily leverage this information. Moreover, the notion of compute units in Kraken fit cleanly with Hadoop-YARN's model in which tasks execute inside containers, an abstraction for a fixed amount of resources. Hadoop also supports the re-spawning of tasks at other Node Managers, a facility already employed by the speculative execution mechanism. Therefore, the Kraken model can be implemented within Hadoop, wherein a container also includes a specification of

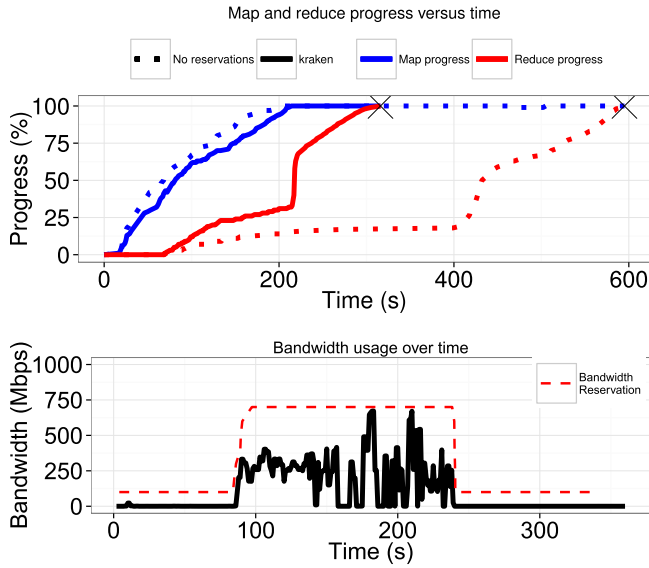


Fig. 11. *Top*: Map and reduce progress of a TeraSort job in best effort conditions (dotted) and when using Kraken (solid lines). The X marks indicate the point of job completion. *Bottom*: Bandwidth utilization for a TeraSort job (black) with online bandwidth reservations (red) via Kraken.

bandwidth, and the Application Master may migrate tasks to different containers in order to satisfy upgrade and/or placement requests.

While we defer a full-fledged implementation of Kraken as future work, in the following, we report on a preliminary prototype, essentially a simplified version of the bandwidth-reservation mechanism demonstrated in PANE [11]: however, while PANE enforces bandwidth reservations as specified by the administrator, we envision Kraken to feed a system like PANE with the exact guarantees to enforce.

We implemented a simple controller that runs inside a virtual machine and uses the Linux `tc` utility in order to make bandwidth reservations. We instrumented the Hadoop source code such that tasks inform the controller prior to executing a shuffle. If there is spare bandwidth to be allocated, the controller increases the corresponding endpoint’s bandwidth reservation. Once the shuffle is completed, the Hadoop framework informs the controller of the same in order to release its reservations. Migrations are not supported yet.

Using this framework, we then executed the TeraSort benchmark against a dataset size of 100 million 100-byte records (a total of  $\sim 10$  GB of data). We co-locate the Hadoop data nodes against a set of VMs that generate UDP flows on the network using `iperf`. The UDP flows are set to last 400 seconds at a throughput of 600 Mbps, which significantly stresses the underlying 1 Gbps network. Hadoop is initially allocated only 100 Mbps of bandwidth, but requests additional bandwidth when shuffles are to be executed (with the UDP tenant given best-effort service, and left to consume the remaining bandwidth throughout). Figure 11 (*top*) indicates the map and reduce progress with and without online reservations. Without online reservations, the UDP tenant interferes heavily with Hadoop’s network usage, thus prolonging the TeraSort job until the UDP flows terminate. With online reservations however, Hadoop requests bandwidth when it needs it, leading

to just under 300 seconds of improvement in job completion time. Figure 11 (*bottom*) indicates the the bandwidth reservations over time as requested by Hadoop (red) and the actual bandwidth used by Hadoop (black). Note that this is in contrast to having to provision for the peak utilization for the entire duration of the run, which would in turn affect efficiency.

## VII. RELATED WORK

**Cloud Network Performance.** Over the last few years, researchers and practitioners have recognized the importance of predictable network performance in a multi-tenant datacenter [26]. One solution is the one adopted by Amazon’s Compute Cluster approach, which avoids multi-tenancy entirely but comes at the cost of reduced efficiency (limited or no resource sharing). Other solutions, instead, extend the max-min *per-flow* based fairness model provided by TCP to support new *per-tenant* fairness models [21], [23], [28], [32]. With these proposals, however, the performance for a given tenant is still dependent on the number of other tenants and their workload, and, hence, cannot be accurately predicted.

In contrast, solutions based on explicit bandwidth reservations [5], [9], [17], [29], [30] allow the tenant to specify the desired bandwidth, usually assuming a Hose model [13], and sometimes even with work-conservation guarantees [18]. The main motivation for our work is that *none* of these solutions allows tenants to update their bandwidth reservation at runtime.

Proteus [9] introduces the concept of *Time-Interleaved Virtual Clusters* which model the time-varying nature of networking requirement of cloud applications such as Hadoop. However, in contrast to Kraken, the number of virtual machines as well as their location in the substrate is *constant* (i.e., fixed) during the entire execution; only the bandwidth reservation between the CUs can be changed over time. Also, it may not be possible to satisfy certain time-varying requests upfront, as they can only be realized *with migrations*. Finally, Proteus is an *offline* approach: In Proteus, each application is profiled first, and the inferred execution patterns are then taken into account when embedding the virtual networks. As we have argued, we believe that this approach is problematic in multi-tenant datacenters, where unexpected events and stragglers are a reality; moreover, it limits the approach to batch-processing type of applications with limited runtimes only (long-running services like interactive web-sites or data stores cannot be modeled). Finally, even in the absence of failures or stragglers, we argue that execution patterns may significantly differ from execution to execution, due to factors such as varying data inputs and differences in data-locality caused by the application.

**Elastic Computing:** A recent class of systems exploit the so-called *time malleability* of many batch processing framework to reconfigure at runtime the amount of resources allocated to each job. For example, Amoeba [2] and Natjam [8] use task preemption to re-distribute at runtime the resources allocated to running jobs. This can be used, for instance, to compensate for deviations from the expected performance due to stragglers [12] or to enable flexible pricing mechanisms [25].

Due to the lack of systems that enable dynamic reconfiguration of network resources, all these systems only focus on computation resources. However, in real workloads, a large fraction of the time of a job is spent on network transfers [24]. Therefore, by disregarding network resources, the effectiveness of these systems is greatly reduced.

We believe that Kraken can be successfully integrated with these systems to provide comprehensive solutions that take into account both compute and network resources. We leave the exploration of these opportunities to future work.

Kraken can also be applied to systems such as Bazaar [19] that provide a job-centric interface and allow the provider to select the best combination of CUs and network resources. By adding the ability of reallocating CUs and network resources at runtime, we can expand the range of scheduling opportunities.

**Embedding:** Many existing systems try to maximize the number of virtual networks that can be hosted concurrently on a given physical infrastructure, while providing the specified resource isolation [6]. Accordingly, embedding and/or scheduling algorithms have been proposed to multiplex virtual networks. Since the underlying problems are often computationally hard in many models [7], approximation algorithms or heuristics are typically used.

This paper focuses on the embedding of virtual clusters in fat-tree networks, and hence, from an embedding algorithm point of view, Oktopus and Proteus are the papers most related to ours. Both systems are based on a collocation strategy which try to place CUs of the same virtual network close to each other, in order to minimize the overall bandwidth allocation. As a side contribution of our current paper, we also present a linear-time algorithm to solve the virtual cluster embedding problem in the fat-tree optimally: the problem considered in [5] and [9] is *not* NP-hard and our algorithm is significantly faster than the algorithms presented in [14] and [31] for other topologies.

**Bibliographic Note.** A shorter version of this paper (without detailed analysis and prototype) appeared at the IEEE INFOCOM 2016 conference [15].

### VIII. DISCUSSION

This paper presented the Kraken system which allows to dynamically scale up and down the bandwidth and compute resources allocated to a cloud application at runtime. Thus, Kraken overcomes the weaknesses of existing solutions, in which resource reservations either cannot be changed [5], [17], [30], in which the entire resource schedule has to be computed at job submission time [9], or in which either only the bandwidth or the compute resources can be adapted, but not both [9], [27], [33].

We described algorithms to find a configurable and optimal trade-off between embedding and reconfiguration costs, and complemented the formal guarantees by simulation and through a preliminary Hadoop prototype.

While we have motivated our approach for batch-processing applications such as MapReduce, the problem is relevant more generally. We also believe that our perspective nicely complements the recent work on *time malleable* systems like

Amoeba [2] and Natjam [8] or scheduling frameworks such as Jokey [12]. Kraken can also be applied to systems such as Bazaar [19] that provide a job-centric interface and allow the provider to select the best combination of CUs and network resources. The ability of reallocating CUs and network resources at runtime can expand the range of scheduling opportunities.

We believe that our work opens several interesting directions for future research. On the theoretical side, it will be interesting to study how to generalize our algorithms beyond fat-tree networks: related work (such as [31]) on single request embeddings without support for migration suggests that polynomial-time algorithms may still exist. The main open question however concerns the study of scheduling algorithms that leverage the Kraken interface to better *schedule* executions over time, also leveraging possible prediction models.

### REFERENCES

- [1] M. Al-Fares, A. Loukissas, and A. Vahdat, "A scalable, commodity data center network architecture," in *Proc. ACM SIGCOMM*, 2008, pp. 63–74.
- [2] G. Ananthanarayanan, C. Douglas, R. Ramakrishnan, S. Rao, and I. Stoica, "True elasticity in multi-tenant data-intensive compute clusters," in *Proc. ACM SOCC*, 2012, Art. no. 24.
- [3] G. Ananthanarayanan, A. Ghodsi, S. Shenker, and I. Stoica, "Effective straggler mitigation: Attack of the clones," in *Proc. USENIX NSDI*, 2013, pp. 185–198.
- [4] G. Ananthanarayanan *et al.*, "Reining in the outliers in Map-Reduce clusters using mantri," in *Proc. USENIX OSDI*, 2010, p. 24.
- [5] H. Ballani, P. Costa, T. Karagiannis, and A. Rowstron, "Towards predictable datacenter networks," in *Proc. ACM SIGCOMM*, 2011, pp. 242–253.
- [6] A. Belbekkouche, M. M. Hasan, and A. Karmouch, "Resource discovery and allocation in network virtualization," *IEEE Commun. Surveys Tuts.*, vol. 14, no. 4, pp. 1114–1128, 4th Quart., 2012.
- [7] C. Chekuri, F. B. Shepherd, G. Oriolo, and M. G. Scutella, "Hardness of robust network design," *Networks*, vol. 50, no. 1, pp. 50–54, 2007.
- [8] B. Cho *et al.*, "Natjam: Design and evaluation of eviction policies for supporting priorities and deadlines in MapReduce clusters," in *Proc. ACM SOCC*, 2013, Art. no. 4.
- [9] D. Xie, N. Ding, Y. C. Hu, and R. Kompella, "The only constant is change: Incorporating time-varying network reservations in data centers," in *Proc. ACM SIGCOMM*, 2012, pp. 199–210.
- [10] G. Even, M. Medina, G. Schaffrath, and S. Schmid, "Competitive and deterministic embeddings of virtual networks," *Theor. Comput. Sci.*, vol. 496, pp. 184–194, Jul. 2013.
- [11] A. Ferguson, A. Guha, C. Liang, R. Fonseca, and S. Krishnamurthi, "Participatory networking: An API for application control of SDNs," in *Proc. ACM SIGCOMM*, Hong Kong, Aug. 2013, pp. 327–338.
- [12] A. D. Ferguson, P. Bodik, S. Kandula, E. Boutin, and R. Fonseca, "Jockey: Guaranteed job latency in data parallel clusters," in *Proc. 7th ACM EuroSys*, 2012, pp. 99–112.
- [13] J. A. Fingerhut, S. Suri, and J. S. Turner, "Designing least-cost non-blocking broadband networks," *J. Algorithms*, vol. 24, pp. 287–309, Aug. 1997.
- [14] C. Fuerst, M. Pacut, P. Costa, and S. Schmid, "How hard can it be?: Understanding the complexity of replica aware virtual cluster embeddings," in *Proc. 23rd IEEE Int. Conf. Netw. Protocols (ICNP)*, Nov. 2015, pp. 11–21.
- [15] C. Fuerst, S. Schmid, L. Suresh, and P. Costa, "Kraken: Online and elastic resource reservations for multi-tenant datacenters," in *Proc. 35th IEEE Conf. Comput. Commun. (INFOCOM)*, Apr. 2016, pp. 1–9.
- [16] A. Greenberg *et al.*, "VL2: A scalable and flexible data center network," in *Proc. ACM SIGCOMM*, 2009, pp. 51–62.
- [17] C. Guo *et al.*, "SecondNet: A data center network virtualization architecture with bandwidth guarantees," in *Proc. ACM CoNEXT*, 2010, Art. no. 15.
- [18] S. Hu *et al.*, "Providing bandwidth guarantees, work conservation and low latency simultaneously in the cloud," in *Proc. IEEE INFOCOM*, Apr. 2016, pp. 1–9.

- [19] V. Jalaparti, H. Ballani, P. Costa, T. Karagiannis, and A. Rowstron, "Bridging the tenant-provider gap in cloud services," in *Proc. SOCC*, 2012, Art. no. 10.
- [20] Y. Kwon, M. Balazinska, B. Howe, and J. Rolia, "SkewTune in action: Mitigating skew in MapReduce applications," *Proc. VLDB Endowment*, vol. 5, no. 12, pp. 1934–1937, 2012.
- [21] L. Popa *et al.*, "FairCloud: Sharing the network in cloud computing," in *Proc. ACM SIGCOMM*, 2012, pp. 187–198.
- [22] K. LaCurtis, J. C. Mogul, H. Balakrishnan, and Y. Turner, "Cicada: Introducing predictive guarantees for cloud networks," in *Proc. 6th USENIX Workshop Hot Topics Cloud Comput. (HotCloud)*, Philadelphia, PA, USA, Jun. 2014, pp. 14–19.
- [23] T. Lam, S. Radhakrishnan, A. Vahdat, and G. Varghese, "NetShare: Virtualizing data center networks across services," Univ. California, San Diego, CA, USA, Tech. Rep. CS2010-0957, 2010.
- [24] M. Chowdhury, M. Zaharia, J. Ma, M. I. Jordan, and I. Stoica, "Managing data transfers in computer clusters with orchestra," in *Proc. ACM SIGCOMM*, 2011, pp. 98–109.
- [25] L. Mai, E. Kalyvianaki, and P. Costa, "Exploiting time-malleability in cloud-based batch processing systems," presented at the LADIS, 2013.
- [26] J. C. Mogul and L. Popa, "What we talk about when we talk about cloud network performance," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 42, no. 5, pp. 44–48, 2012.
- [27] R. Nathuji, A. Kansal, and A. Ghaffarkhah, "Q-clouds: Managing performance interference effects for QoS-aware clouds," in *Proc. EuroSys*, 2010, pp. 237–250.
- [28] L. Popa *et al.*, "ElasticSwitch: Practical work-conserving bandwidth guarantees for cloud computing," in *Proc. ACM SIGCOMM*, 2013, pp. 351–362.
- [29] B. Raghavan, K. Vishwanath, S. Ramabhadran, K. Yocum, and A. C. Snoeren, "Cloud control with distributed rate limiting," in *Proc. ACM SIGCOMM*, 2007, pp. 337–348.
- [30] H. Rodrigues, J. R. Santos, Y. Turner, P. Soares, and D. O. Guedes, "Gatekeeper: Supporting bandwidth guarantees for multi-tenant data-center networks," in *Proc. 3rd Conf. I/O Virtualization (WIOV)*, 2011, p. 6.
- [31] M. Rost, C. Fuerst, and S. Schmid, "Beyond the stars: Revisiting virtual cluster embeddings," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 45, no. 3, pp. 12–18, 2015.
- [32] A. Shieh, S. Kandula, A. Greenberg, and C. Kim, "Seawall: Performance isolation for cloud datacenter networks," in *Proc. USENIX HotCloud*, 2010, pp. 1–7.
- [33] K. C. Webb, A. Roy, K. Yocum, and A. C. Snoeren, "Blender: Upgrading tenant-based data center networking," in *Proc. ANCS*, 2014, pp. 65–75.
- [34] M. Zaharia, A. Konwinski, A. D. Joseph, R. Katz, and I. Stoica, "Improving MapReduce performance in heterogeneous environments," in *Proc. USENIX OSDI*, 2008, pp. 29–42.
- [35] *Measuring EC2 System Performance*. Accessed: Aug. 1, 2017. [Online]. Available: <http://goo.gl/V5zhEd>



**Carlo Fuerst** received the Diploma and Ph.D. degrees in computer science from the Technical University of Berlin, Germany, in 2011 and 2016, respectively. He is currently with the SAP Innovation Center Potsdam. His research interests revolve around the reliable performance for elastic cloud environments and algorithms for virtual network embeddings.



**Stefan Schmid** received the M.Sc. and Ph.D. degrees from ETH Zurich, Switzerland, in 2004 and 2008, respectively. In 2009, he held a post-doctoral position at the Technical University of Munich and the University of Paderborn. From 2009 to 2015, he was a Senior Research Scientist with Telekom Innovations Laboratories, Berlin, Germany. From 2015 to 2017, he was an Associate Professor with Aalborg University, Denmark. He is currently a Professor with the University of Vienna, Austria. His research interests revolve around the fundamental and algorithmic problems of networked and distributed systems.



**Lalith Suresh** received the joint master's degree in distributed computing from KTH, Stockholm, and IST, Lisbon, and the Ph.D. degree from the Technical University of Berlin, under the supervision of Prof. A. Feldmann and Prof. M. Canini. He is currently a Research Scientist with the VMware Research Group, Palo Alto.



**Paolo Costa** received the M.Sc. and Ph.D. degrees in computer engineering from the Politecnico di Milano. He has been a Research Faculty with the Imperial College London and a Post-Doctoral Researcher with the Computer Systems Group, Vrije Universiteit Amsterdam. He is currently a Researcher with the Systems and Networking Group, Microsoft Research, and an Honorary Lecturer with the Department of Computing, Imperial College London. His research interests include the intersection of distributed systems and networking with particular emphasis on large-scale data centers.