

RESOURCE MANAGEMENT IN BIG DATA PROCESSING SYSTEMS

S. Tang, B. He, H. Liu, B.-S. Lee

7.1 INTRODUCTION

In many application domains such as social networks and bioinformatics, data is being gathered at unprecedented scale. Efficient processing for Big Data analysis poses new challenges for almost all aspects of state-of-the-art data processing and management systems. For example, there are a few challenges as follows: (i) the data can be arbitrarily complex structures (eg, graph data) and cannot be efficiently stored in relational database; (ii) the data access of large-scale data processing are frequent and complex, resulting in inefficient disk I/O accesses or network communications; and (iii) last but not least, a variety of unpredictable failure problems must be tackled in the distributed environment, so the data processing system must have a fault tolerance mechanism to recovery the task computation automatically.

Cloud computing has emerged as an appealing paradigm for Big Data processing over the Internet due to its cost effectiveness and powerful computational capacity. Current infrastructure-as-a-service (IaaS) clouds allow tenants to acquire and release resources in the form of virtual machines (VMs) on a pay-as-you-use basis. Most IaaS cloud providers such as Amazon EC2 offer a number of VM types (such as *small*, *medium*, *large*, and *extra-large*) with fixed amount of CPU, main memory, and disk. Tenants can only purchase fixed-size VMs and increase/decrease the number of VMs when the resource demands change. This is known as *T-shirt and scale-out model* [1]. However, the T-shirt model leads to inefficient allocation of cloud resources, which translates to higher capital expenses and operating costs for cloud providers, as well as an increase of monetary cost for tenants. First, the granularity of resource acquisition/release is coarse in the sense that the fix-sized VMs are not tailored for cloud applications with dynamic demands delicately. As a result, tenants need to overprovision resource (costly), or risk performance penalty and service level agreement (SLA) violations. Second, elastic resource scaling in clouds [2], also known as a scale-out model, is also costly due to the latencies involved in VM instantiating [3] and software runtime overhead [4]. These costs are ultimately borne by tenants in terms of monetary cost or performance penalty.

Resource sharing is a classic and effective approach to resource efficiency. As more and more Big Data applications with diversifying and heterogeneous resource requirements tend to deploy in the cloud [5,6], there are vast opportunities for resource sharing [1,7]. Recent work has shown that fine-grained and dynamic resource allocation techniques (eg, resource multiplexing or overcommitting [1,8–11]) can significantly improve the resource utilization compared to T-shirt model [1]. As adding/removing resources is directly performed on the existing VMs, a fine-grained resource allocation is also known as a scale-up model, and the cost tends to much smaller compared to the scale-out model. Unfortunately, current IaaS clouds do not offer resource sharing among VMs, even if those VMs belong

to the same tenant. Resource sharing models are needed for better cost efficiency of tenants and higher resource utilization of cloud providers.

Researchers have been actively proposing many innovative solutions to address the new challenges of large-scale data processing and resource management. In particular, a notable number of large-scale data processing and resource management systems have recently been proposed. The aims of this chapter are (i) to introduce canonical examples of large data processing, (ii) to make an overview of existing data processing and resource management systems/platforms, and more importantly, (iii) to make a study on the economic fairness for large-scale resource management on the cloud, which bridges large-scale resource management and cloud-based platforms. In particular, we present some desirable properties including sharing incentive, truthfulness, resource-as-you-pay fairness, and pareto efficiency to guide the design of fair policy for the cloud environment.

The chapter is organized as follows: We first present several types of resource management for Big Data processing in [Section 7.2](#). In [Section 7.3](#), we list some representative Big Data processing platforms. [Section 7.4](#) presents the single resource management in the cloud, following by [Section 7.5](#) that introduces multiresource management in the cloud. For the completeness of discussions, [Section 7.6](#) gives an introduction to existing work on resource management. A discussion of open problems is provided in [Section 7.7](#). We conclude the chapter in [Section 7.8](#).

7.2 TYPES OF RESOURCE MANAGEMENT

Resource management is a general and fundamental issue in computing systems. In this section, we present the resource management for typical resources including CPU, memory, storage, and network.

7.2.1 CPU AND MEMORY RESOURCE MANAGEMENT

Current supercomputers and data centers (eg, Amazon EC2) generally consist of thousands of computing machines. At any time, there are tens of thousands of users running their high-performance computing applications (eg, MapReduce [12], Message Passing Interface (MPI), Spark [13]) on it. The efficient resource management of the computing resources such as CPU, memory is nontrivial for high performance and fairness. Typically, the resource management includes resource discovery, resource scheduling, resource allocation, and resource monitoring. Resource discovery identifies the suitable computing resources in which machines that match the user's request. Resource scheduling selects the best resource from the matched computing resources. It actually identifies the physical resource where the machines are to be created to provision the resources. Resource allocation allocates the selected resource to the job or task of the user's request; it means the job submission to the selected cloud resource. After the submission of the job, the resource is monitored.

There are a number of resource management tools available for supercomputing. For example, Simple Linux Utility for Resource Management (SLURM) is a highly scalable resource manager widely used in supercomputers [14]. It allocates exclusive and/or nonexclusive access to resources (computer nodes) to users for some duration of time so they can perform work. Second, it provides a framework for starting, executing, and monitoring work (typically a parallel job) on a set of allocated nodes. Finally, it arbitrates contention for resources by managing a queue of pending work. For data-intensive computing in data center, YARN [15] and Mesos [16] are two popular resource management systems.

7.2.2 STORAGE RESOURCE MANAGEMENT

Storage resource management (SRM) is a proactive approach to optimizing the efficiency and speed with which available drive space is utilized in a storage area network, which is a dedicated high-speed network (or subnetwork) that interconnects and presents shared pools of storage devices to multiple servers [17]. The SRM software can help a storage administrator automate data backup, data recovery, and Storage Area Network (SAN) performance analysis. It can also help the administrator with configuration management and performance monitoring, forecast future storage needs more accurately, and understand where and how to use tiered storage, storage pools, and thin provisioning.

7.2.3 NETWORK RESOURCE MANAGEMENT

Managing and allocating the network flows to different applications/users is nontrivial work. Particularly, software-defined networking [18] is a popular approach that allows network administrators to manage network services through the abstraction of lower-level functionality. This is done by decoupling the system that makes decisions about where traffic is sent (the control plane) from the underlying systems that forward traffic to the selected destination (the data plane).

7.3 BIG DATA PROCESSING SYSTEMS AND PLATFORMS

In the era of Big Data, characterized by the unprecedented volume of data, the velocity of data generation, and the variety of the structure of data, support for large-scale data analytics constitutes a particularly challenging task. To address the scalability requirements of today's data analytics, parallel shared-nothing architectures of commodity machines (often consisting of thousands of nodes) have been lately established as the de facto solution. Various systems have been developed mainly by the industry to support Big Data analysis, including MapReduce [12], Pregel [19], Spark [13], etc. In this section, we give a brief survey of some representative solutions.

7.3.1 HADOOP

Hadoop [20] is an open source Java implementation of MapReduce [12], which is a popular programming model for large-scale data processing proposed by Google. It runs on top of a distributed file system called HDFS, which splits input data into multiple blocks of fixed size (typically 64 MB) and replicates each data block several times across computing nodes. Users can submit MapReduce jobs to the Hadoop cluster. The Hadoop system breaks each job into multiple map tasks and reduce tasks, with its map tasks computed before its reduce tasks. Each map task processes (ie, scans and records) a data block and produces intermediate results in the form of key-value pairs.

Moreover, Hadoop has evolved to the next generation of Hadoop (ie, Hadoop MRv2) called YARN [15], as a large-scale data operating platform and cluster resource management system. There is a new architecture for YARN, which separates the resource management from the computation model. Such a separation enables YARN to support a number of diverse data-intensive computing frameworks including Dryad [21], Giraph [22], Spark [13], Storm [23], and Tez [24]. In YARN's architecture, there is a global master named ResourceManager (RM) and a set of per-node slaves

called NodeManagers (NM), which forms a generic system for managing applications in a distributed manner. The RM is responsible for tracking and arbitrating resources among applications. In contrast, the NM has responsibility for launching tasks and monitoring the resource usage per slave node. Moreover, there is another component called ApplicationMaster, which is a framework-specific entity. It is responsible for negotiating resources from the RM and working with the NM to execute and monitor the progress of tasks. Particularly, all resources of YARN are requested in the form of container, which is a logical bundle of resources (eg, 1 CPUs, 2G memory). As a multitenant platform, YARN organizes user submitted applications into queues and share resources between these queues. Users can set their own queues in a configuration file provided by YARN.

7.3.2 DRYAD

Dryad [21] is a distributed execution engine that simplifies the process of implementing data-parallel applications to run on a cluster. The original motivation for Dryad was to execute data mining operations efficiently, which has also lead to technologies such as MapReduce or Hadoop. Dryad is a general-purpose execution engine and can also be used to implement a wide range of other application types, including time series analysis, image processing, and a variety of scientific computations. The computation is structured as a directed graph: Programs are graph vertices, while the channels are graph edges. A Dryad job is a graph generator, which can synthesize any directed acyclic graph. These graphs can even change during execution in response to important events in the computation. Dryad handles job creation and management, resource management, job monitoring and visualization, fault tolerance, reexecution, scheduling, and accounting.

7.3.3 PREGEL

Pregel [19] is a specialized model for iterative graph applications. In Pregel, a program runs as a series of coordinated supersteps. With each superstep, each vertex in the graph runs a user function that can update state associated with the vertex, change the graph topology, and send messages to other vertices for use in the next superstep. This model can express many graph algorithms, including shortest paths, bipartite matching, and PageRank.

7.3.4 STORM

Storm [23] is a distributed real-time computation system. In a similar way to how Hadoop provides a set of general primitives for doing batch processing, Storm provides a set of general primitives for doing real-time computation, which greatly eases the writing of parallel real-time computation. It can be used for processing messages and updating databases (stream processing), performing a continuous query on data streams and streaming the results into clients (continuous computation), parallelizing an intense query like a search query on the fly (distributed Remote Procedure Calls (RPC)), and more. Storm's small set of primitives satisfies a stunning number of user needs. For example, Storm scales to massive numbers of messages per second. To scale a topology, all you have to do is add machines and increase the parallelism settings of the topology. As an example of Storm's scale, one of Storm's initial applications processed 1,000,000 messages per second on a 10-node cluster, including hundreds of database calls per second as part of the topology.

7.3.5 SPARK

Spark [13] is a fast in-memory data processing system that achieves high performance for applications through caching data in memory (or disk) for data sharing across computation stages. It is achieved with the *resilient distributed dataset* (RDD) in-memory storage abstraction for computing data, which is a read-only, partitioned collection of records [25]. Each RDD is generated from data in stable storage or other RDDs through *transformation* operations such as *map*, *filter*, and *reduceByKey*. Notably, the *transformation* is a *lazy* operation that only defines a new RDD without immediately computing it. To launch RDD computation, Spark provides another set of *action* operations such as *count*, *collect*, and *save*, which return a value to the application or export data to a storage system. For RDD caching, Spark offers five storage levels, ie, *memory only*, *memory and disk*, *memory only ser*, *memory and disk ser*, and *disk only*. For a Spark application in execution, the Spark system will spawn a master called *driver*, which is responsible for defining and managing RDDs, and a set of slavers called *executors*, which perform the computations dynamically. The Spark applications can run on either YARN, Mesos, local, or standalone cluster modes. In this paper, we focus on the standalone cluster mode.

7.3.6 SUMMARY

In this section, we summarize some representative Big Data processing systems elaborated in the previous section in terms of computation model, in-memory computation, resource management type (ie, CPU, memory, storage, and network), and fairness in Table 1.

There are some other high-level as well as application-specific systems that are built on top of previous data computing systems to form an ecosystem for a variety of applications. For example, for Hadoop, Apache Pig [26], and Hive [27] are both Structured Query Language (SQL)-like systems that are running on it to support analytical data querying processing. HBase [28] is a NoSql database system built on top of Hadoop system. Apache Giraph [22] is an iterative graph processing system running on Hadoop. Similarly, for Spark, Shark [29], and Spark SQL [30] are two analytical data query systems built on Spark, while Graphx [31] is a graph processing system for graph applications. We have also witnessed some other data processing systems/platforms that are running on currently emerging computing devices, such as GPUs. For example, as an extension of Pregel for GPU platform, a general-purpose programming framework called Medusa [32] has been developed for graph applications.

Table 1 Comparison of Representative Big Data Processing Systems

Systems	Computation Model	In-Memory Computation	Resource Management Type				Resource Fairness	
			CPU	Memory	Storage	Network	Single	Multiple
Hadoop	MapReduce	No	Yes	Yes	No	No	Yes	Yes
Dryad	Dryad	No	Yes	No	No	No	Yes	No
Pregel	Pregel	No	Yes	No	No	No	Yes	No
Storm	Storm	No	Yes	No	No	No	Yes	No
Spark	RDD	Yes	Yes	Yes	No	No	Yes	Yes

7.4 SINGLE-RESOURCE MANAGEMENT IN THE CLOUD

Single-resource management refers to the management of single-resource type. This is the basic form of resource management. For example, Hadoop manages CPU resources of a computing cluster in the form of slots.

Cloud computing has emerged as a popular platform for users to compute their large-scale data applications, attracting from its merits such as flexibility, elasticity, and cost efficiency. At any time, there can be tens of thousands of users concurrently running their large-scale, data-intensive applications on the cloud. Users pay money on the basis of their resource usage. To meet different users' needs, cloud providers generally offer several options of price plans (eg, on-demand price, reserved price). When users have a short-term computation requirement (eg, 1 week), they can choose an on-demand price plan that charges compute resources by each time unit (eg, an hour) with a fixed price. In contrast, if users have long-term computation requests (eg, 1 year), choosing a reserved price plan can enable them to have a significant discount on the hourly charge for the resources in comparison to the on-demand one, thereby saving money.

To improve the resource utilization and in turn the cost efficiency, resource sharing is an effective approach [7]. Consider the reserved resources for example: With a reservation plan, users need to pay a one-time fee for a long time (eg, 1 or 3 years), and in turn get a significant discount on the hourly usage charge. However, to achieve the full cost savings, users must commit to having a high utilization. In practice, the resource demand of a user can fluctuate over time, and the resources cannot be fully utilized all the time from the perspective of an individual user. With resource sharing, users can complement each other in the resource usage in a shared cluster and the resource utilization problem can be thereby solved. To make resource sharing possible among users, the resource allocation fairness is a key issue.

As shown in Table 1, there are single- and multiple-resource fairness, both of which are supported by some of computing systems. In this section, we mainly focus on the fairness of single-resource management and defer the fairness for multiple resource management to Section 7.5. Notably, we observe that the fair policies implemented in these systems are all memoryless (ie, allocating resources fairly at instant time without considering history information). We refer to those schedulers as *memoryless resource fairness* (MLRF). In the following subsections, we first present several desirable resource allocation properties for cloud computing. Next, we show the problems for MLRF. Finally, we explore a new policy to address it.

7.4.1 DESIRED RESOURCE ALLOCATION PROPERTIES

This section presents a set of desirable properties that we believe any cloud-oriented resource allocation policy in a shared pay-as-you-use system should meet

- *Sharing incentive*: Each client should be better off sharing the resources via group-buying with others, than exclusively buying and using the resources individually. Consider a shared pay-as-you-use computing system with n clients over t period time. Then a client should not be able to get more than t/n resources in a system partition consisting of $1/n$ of all resources.
- *Nontrivial workload incentive*: Clients should reap benefits by submitting nontrivial workloads and yielding unused resources to others when not needed. Otherwise, they may unnecessarily hold unneeded resources under their share by running trivial tasks in a shared computing environment.

- *Resource-as-you-pay fairness*: The resource that clients gain should be proportional to their payment. This property is important, as it is a resource guarantee to clients.
- *Strategy-proofness*: Clients should not be able to get benefits by lying about their resource demands. This property is compatible with sharing incentive and resource-as-you-pay fairness, since no client can obtain more resources by lying.
- *Pareto efficiency*: In a shared resource environment, it is impossible for a client to get more resources without decreasing the resource of at least one client. This property can ensure the system resource utilization to be maximized.

7.4.2 PROBLEMS FOR EXISTING FAIRNESS POLICIES

One of the most popular fair allocation policies is (*weighted*) *max-min fairness* [7], which maximizes the minimum resource allocation obtained by a user in a shared computing system. It has been widely used in many popular large-scale data processing frameworks such as Hadoop [20], YARN [15], Mesos [16], and Dryad [21]. Unfortunately, we observe that the fair policies implemented in these systems are all *memoryless* (ie, allocating resources fairly at instant time without considering history information). We refer those schedulers as MLRF. MLRF is *not* suitable for a cloud computing system due to the following reasons.

Trivial workload problem

In a pay-as-you-use computing system, we should have a policy to incentivize group members to submit *nontrivial* workloads that they really need (see *Nontrivial workload incentive* property in Section 7.4.1). For MLRF, there is an implicit assumption that all users are unselfish and honest towards their requested resource demands, which, however, is often not true in real world. It can cause trivial workload problems with MLRF. Consider two users A and B sharing a system. Let D_A and D_B be the *true* workload demand for A and B at time t_0 , respectively. Assume that D_A is less than its share,¹ while D_B is larger than its share. In that case, it is possible that A is selfish and will try to possess all of his/her share by running some trivial tasks (eg, running some duplicated tasks of the experimental workloads for double checking) so that his/her extra unused share will not be preempted by B, causing the inefficiency problem of running nontrivial workloads and also breaking the sharing incentive property (see the definition in Section 7.4.1).

Strategy-proofness problem

It is important for a shared system to have a policy to ensure that no group member can get any benefits by lying (see *Strategy-proofness* in Section 7.4.1). We argue that MLRF cannot satisfy this property. Consider a system consisting of three users A, B, and C. Assume A and C are honest whereas B is not. It could happen at a time that the *true* demands of both A and B are fewer than their own shares, while C's *true* demands exceed his/her share. In that case, A yields his/her unused resources to others honestly. But B will provide *false* information about his/her demand (eg, far larger than his/her share) and compete with C for unused resources from A. Lying benefits B, hence violating strategy-proofness. Moreover, it will break the sharing incentive property if all other users also lie.

¹By default, we refer to the *current* share at the designated time (eg, t_0), rather than the *total* share accumulated over time.

Resource-as-you-pay unfairness problem

For group-buying resources, we should ensure that the total resource received by each member is proportional to his/her monetary cost (see *Resources-as-you-pay fairness* in Section 7.4.1). Due to the varied resource demands (eg, workflows) for a user at different time, MLRF cannot achieve this property. Consider two users A and B: At time t_0 , it could happen that the demand D_A is less than its share; hence, its extra unused resource will be possessed by B (ie, lent to B) according to the work-conserving property of MLRF. Next at time t_1 , assume that A's demand D_A becomes larger than its share. With MLRF, user A can only use her current share (ie, cannot get lent resources at t_0 back from B), if D_B is larger than its share, due to *memoryless*. If this scenario often occurs, it will be unfair for A to get the amount of resources that she should have obtained from a long-term view.

7.4.3 LONG-TERM RESOURCE ALLOCATION POLICY

In this section, we first give a motivation example to show that MLRF is *not* a suitable cloud computing system. Then we propose long-term resource fairness (LTRF), a cloud-oriented allocation policy to address the limitations of MLRF and meet the desired properties described in Section 7.4.1.

Motivation example

Consider a shared computing system consisting of 100 resources (eg, 100GB RAM) and two users A and B with equal share of 50GB each. As illustrated in Table 2, assume that the new requested demands at time t_1 , t_2 , t_3 and t_4 for client A are 20, 40, 80 and 60, and for client B are 100, 60, 50 and 50, respectively. With MLRF, we see in Table 2A that, at t_1 , the total demand and allocation for A are both 20. It lends 30 unused resources to B and thus 80 allocations for B. The scenario is similar at t_2 . Next at t_3 and t_4 , the total demand for A becomes 80 and 90, bigger than its share of 50. However, it can only get 50 allocations based on MLRF, being *unfair* for A, since the total allocations for A and B become 160(=20+40+50+50) and 240(=80+60+50+50) at time t_4 , respectively. Instead, if we adopt LTRF, as shown in Table 2B, the total allocations for A and B at t_4 will finally be the same (eg, 200), being *fair* for A and B.

LTRF scheduling algorithm

Algorithm 1 shows pseudocode for LTRF scheduling. It considers the fairness of total allocated resources consumed by each client instead of currently allocated resources. The core idea is based on the “loan (lending) agreement” [33] with free interest. That is, a client will yield unused resources to others as a *lend* manner at a time. When the client needs resources at a later time, he/she should get the resources back from others that were yielded before (ie, *return* manner). In our previous two-client example with LTRF in Table 2B, client A first lends his/her unused resources of 30 and 10 to client B at time t_1 and t_2 , respectively. However, at t_3 and t_4 , the client has a large demand and then collects all 40 extra resources back from B that were lent before, making *fair* between A and B.

Due to the *lending agreement* of LTRF, in practice, when A yields the unused resources at t_1 and t_2 , B might not want to possess extra unused resources from A immediately. In that case, the total allocations for A and B will be 160(=20+40+50+50) and 200(=50+50+50+50) at time t_4 , causing the inefficiency problem for the system utilization. To solve this problem, we propose a discount-based approach. The idea is that, anybody possessing extra unused resources from others will have a discount

Table 2 A Comparison Example of *MemoryLess Resource Fairness (MLRF)* and *Long-Term Resource Fairness (LTRF)* in a Shared Computing System Consisting of 100 Computing Resources for Two Users A and B)

	Client A					Client B				
	Demand		Allocation		Preempt	Demand		Allocation		Preempt
	New	Total	Current	Total		New	Total	Current	Total	
(A) Allocation results based on <i>MLRF</i> . <i>Total Demand</i> refers to the sum of the new demand and accumulated remaining demand in previous time										
t_1	20	20	20	20	-30	100	100	80	80	+30
t_2	40	40	40	60	-10	60	80	60	140	+10
t_3	80	80	50	10	0	50	70	0	190	0
t_4	60	90	50	160	0	50	70	50	240	0
(B) Allocation results based on <i>LTRF</i>										
t_1	20	20	20	20	-30	100	100	80	0	+30
t_2	40	40	40	60	-10	60	80	60	140	+10
t_3	80	80	80	140	+30	50	70	20	160	-30
t_4	60	60	60	200	+10	50	100	40	200	-10
(C) Counted allocation results under <i>discount</i> -based approach of <i>LTRF</i> . There is a discount (eg, 50%) for the extra unused resources to incentivize clients to preempt resources actively for system utilization maximization. In this example, although the <i>counted</i> allocations for A and B are 180, their real allocations are both 200, which is the same as Table 2B										
t_1	20	20	20	20	-30	100	100	65	65	+30
t_2	40	40	40	60	-10	60	80	55	120	+10
t_3	80	80	65	125	+30	50	70	20	140	-30
t_4	60	60	55	180	+10	50	100	40	180	-10

(eg, 50%) on resource counting. It will incentivize B to preempt extra unused resources from A, since it is cheaper than his/her own share of resources. From a perspective from A, it also does not get resource lost, as it can get the same discount on future resource use later.

Table 2C demonstrates this point. It shows the discounted resource allocation for each client over time by discounting the possessed extra unused resource. At time t_1 , A yields his/her 30 unused resources to B, and B's discounted resources are 65(=50 + 30 * 50%) instead of 80(=50 + 30). Similarly at t_3 , it preempts 30 resources from B and its discounted resources are 65(50 + 30 * 50%). Still, both of them are *fair* at time t_4 .

ALGORITHM 1 LTRF PSEUDOCODE

- 1: R : total resources available in the system.
- 2: $\vec{R} = (\vec{R}_1, \dots, \vec{R}_n)$: current allocated resources. \vec{R}_i denotes the current allocated resources for client i .

```

3:  $U = (u_1, \dots, u_n)$  : total used resources, initially 0.  $u_i$  denotes the total resource consumed by client  $i$ .
4:  $W = (w_1, \dots, w_n)$  : weighted share.  $w_i$  denotes the weight for client  $i$ .
5: while there are pending tasks do
6:   Choose client  $i$  with the smallest total weighted resources of  $u_i/w_i$ .
7:    $d_i \leftarrow$  the next task resource demand for client  $i$ .
8:   if  $\bar{R} + d_i \leq R$  then
9:      $\bar{R}_i \rightarrow \bar{R}_i + d_i$ . /*Update current allocated resources.*/
10:    Update the total resource usage  $u_i$  for client  $i$ .
11:    Allocate resource to client  $i$ .
12:   else
13:     /*The system is fully utilized.*/
14:     Wait until there is a released resource  $r_i$  from client  $i$ .
15:      $\bar{R}_i \leftarrow \bar{R}_i - r_i$ . /*Update current allocated resources */

```

7.4.4 EXPERIMENTAL EVALUATION

We ran our experiments in a cluster consisting of 10 compute nodes, each with 2 Intel X5675 CPUs (6 CPU cores per CPU with 3.07 GHz), 24 GB DDR3 memory, and 56 GB hard disks. The latest version of YARN-2.2.0 is chosen in our experiment, used with a two-level hierarchy. The first level denotes the root queue (containing one master node and nine slave nodes). For each slave node, we configure its total memory resources with 24 GB. The second level denotes the applications (ie, workloads). We ran a macrobenchmark consisting of four different workloads. Thus, four different queues are configured in YARN/LTYARN, namely, *Facebook*, *Purdue*, *Spark*, *HIVE/TPC-H*, corresponding to the following workloads, respectively. (1) A MapReduce instance with a mix of small and large jobs based on the workload at Facebook. (2) A MapReduce instance running a set of large-sized batch jobs generated with Purdue MapReduce Benchmarks Suite [34]. (3) Hive running a series of TPC-H queries. (4) Spark running a series of machine learning applications.

A good sharing policy should be able to first minimize the sharing loss, then maximize the sharing benefit as much as possible (ie, sharing incentive). We make a comparison between MLRF and LTRF for four workloads over time in Fig. 1. All results are relative to the static partition case (without sharing) with sharing benefit/loss degrees of zero. Fig. 1A and B present the sharing benefit/loss degrees, respectively, for MLRF and LTRF. The following observations can be obtained: first, the sharing policies of both MLRF and LTRF can bring sharing benefits for queues (workloads). This is due to the sharing incentive property, ie, each queue has an opportunity to consume more resources than her share at a time, which is better than running at most all of her shared partition in a nonshared partition system. Second, LTRF has a much better result than MLRF. Specifically, Fig. 1A indicates that the sharing loss problem for MLRF is constantly available until all the workloads complete (eg, about -0.5 on average). In contrast, there is no more sharing loss problem after 650 s for LTRF (ie, all workloads get sharing benefits after that). The major reason is that MLRF does not consider historical resource allocation. In contrast, LTRF is a history-based fairness resource allocation policy. It can dynamically adjust the allocation of resources to each queue in terms of historical consumption and lending agreement so that each queue can obtain the required amount of total resources over time. Finally, regarding the sharing loss problem at the early stage (eg, 0 ~ 650 s) of LTRF in Fig. 1B, it is mainly due to the unavoidable waiting allocation problem at the starting stage (ie, the first workload possess all resources, leading late arriving workloads to wait until some tasks complete and release resources).

The problem exists in both MLRF and LTRF. Still, LTRF can smooth this problem until it disappears over time via a lending agreement, while MLRF cannot.

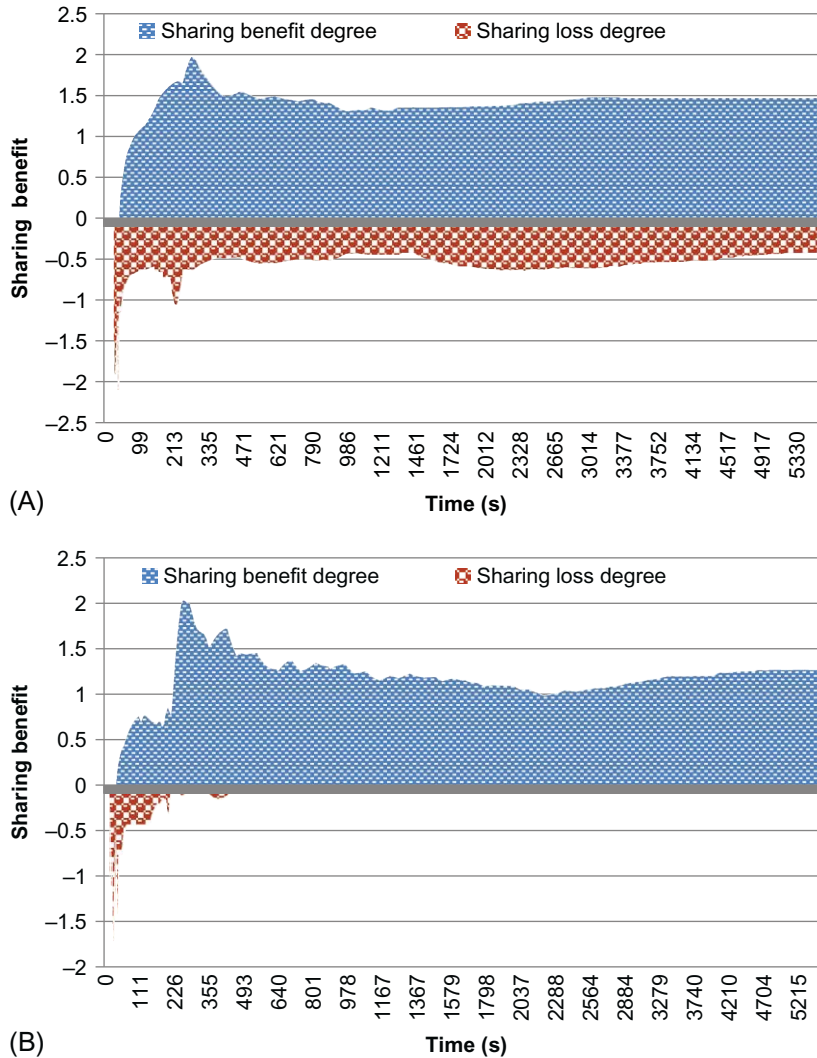


FIG. 1

Comparison of fairness results over time for workloads under MLRF and LTRF in YARN. All results are relative to the static partition scenario (ie, nonshared case) whose sharing benefit/loss is zero. (A) and (B) show the overall benefit/loss relative to the nonsharing scenario.

7.5 MULTIRESOURCE MANAGEMENT IN THE CLOUD

Despite the resource sharing opportunities in the cloud, resource sharing, especially for *multiple* resource types (ie, multiresource), poses several important and challenging problems in pay-as-you-use commercial clouds. (1) *Sharing incentive*. In a shared cloud, a tenant may have concerns about the gain/loss of her asset in terms of resource. (2) *Free riding*. Tenants may deliberately buy fewer resources than

their demand and always expect to benefit from others' contribution (ie, unused resource). However, free riders would seriously hurt other tenants' sharing incentive. (3) *Lying*. When there exists resource contention, a tenant may lie about resource demands for greater benefit. Lying also hurts tenants' sharing incentive. (4) *Gain-as-you-contribute fairness*. It is important to guarantee that the allocations obey the rule "more contribution, more gain." In summary, those problems are eventually attributed to economic fairness of resource sharing in IaaS clouds. Unfortunately, the popular allocation policies such as (weighted) max-min fairness (WMMF) [35] and dominant resource fairness (DRF) [7] cannot address all the problems of resource sharing in IaaS.

This section introduces F2C, a cooperative resource management system for IaaS clouds. F2C exploits statistical resource complementarity to group tenants together in order to realize the resource sharing opportunities, and it adopts a novel resource allocation policy to guarantee fairness of resource sharing. Particularly, we describe *reciprocal resource fairness* (RRF), a generalization of max-min fairness to multiple resource types [36]. The intuition behind RRF is that tenant should preserve their assets while maximizing the resource usage in a cooperative environment. Different types of resource are advocated to trade among different tenants and to share among different VMs belonging to the same tenant. For example, a tenant can trade his/her unused CPU share for another tenant's overprovisioned memory share. In this section, we consider two major kinds of resources, including CPU and main memory. Resource trading (RT) can maximize tenants' benefit from resource sharing.

RRF decomposes the multiresource fair sharing problem into a combination of two complementary mechanisms: intertenant resource trading (IRT) and intratenant weight adjustment (IWA). These mechanisms guarantee that tenants only allocate minimum shares to their nondominant demands and maximize the share allocations on the contended resource. Moreover, RRF is able to achieve some desirable properties of resource sharing, including sharing incentive, gain-as-you-contribute fairness, and strategy-proofness (guarding against free riding and lying).

In the following, we first introduce the resource allocation model in F2C. Second, we analyze the problems of two popular resource allocation policies including WMMF [35] and DRF [7]. Third, we present RRF, the resource allocation model in F2C. Finally, we evaluate F2C and show how RRF can address resource fair allocation problems in IaaS clouds.

7.5.1 RESOURCE ALLOCATION MODEL

We consider the resource-sharing model in multitenant cloud environments, where tenants may rent several VMs to host their applications and VMs have multiresource demands. By *multiresource*, we mean resources of *different resource types*, instead of multiple units of the same resource type. In this paper, we mainly consider two resource types: CPU and main memory. Tenants can have different *weights* (or *shares*) to the resource. The *share* of a tenant reflects the tenant's priority relative to other tenants. A number of tenants can form a resource pool based on the opportunities of resource sharing. The VMs of these tenants then share the same resource pool with negotiated resource *shares*, which are determined by tenants' payments.

In our resource allocation model, each unit of resource is represented by a number of shares. To simplify multiresource allocation, we assume each unit of resource (such as 1 compute unit² or 1 GB

²One EC2 compute unit provides equivalent CPU capacity of a 1.0–1.2GHz 2007 Opteron or 2007 Xeon processor, according to Amazon EC2.

RAM) has its fixed share according to its market price. A study on Amazon EC2 pricing data [37] had indicated that the hourly unit cost for 1 GB memory is twice as expensive as 1 EC2 compute unit. A tenant's *asset* is then defined as the aggregate shares that he/she pays for.

In the following, we use an example to demonstrate our resource allocation model. Fig. 2 shows three tenants colocated on two physical hosts. Each tenant has two VMs. The shares of different resources (eg, CPU, memory) are uniformly normalized based on their market prices [37]. To some extent, a tenant actually purchases resource shares instead of fix-sized resource capacity. A share of a VM reflects the VM's priority relative to other VMs, cloud providers can directly use shares as billing and resource allocation policies. The concrete design of those policies is beyond the scope of this paper. Thus, we simply define a function f_1 to translate tenants' payments into shares $payment \xrightarrow{f_1} share$, and another function f_2 to translate shares into the resource capacity $share \xrightarrow{f_2} resource$. For example, in Fig. 2, 1 compute unit and 1 GB memory are priced at 100 and 200 shares, respectively. If VM1 is initialized with 3 compute units and 2 GB of memory, VM1 is allocated with total $100 \times 3 + 200 \times 2 = 700$ shares.

This model enables fine-grained resource allocation based on shares, and thus provides the opportunities for dynamic resource sharing. Now, the fairness has become a major concern in such a shared system. Informally, we define a kind of *economic fairness*: tenants should try to maximize their aggregate multiresource shares if there are unsatisfied resource demands.

We find that RT between different tenants can reinforce the economic fairness of resource sharing. Normalizing multiple resources with uniform shares provides advantages to facilitate RT. For example, a tenant can trade overprovisioned CPU shares for other tenants' underutilized memory shares. In Fig. 2, VM1 may trade its 200 CPU shares for VM3's 100 memory shares. RT can prevent tenants from

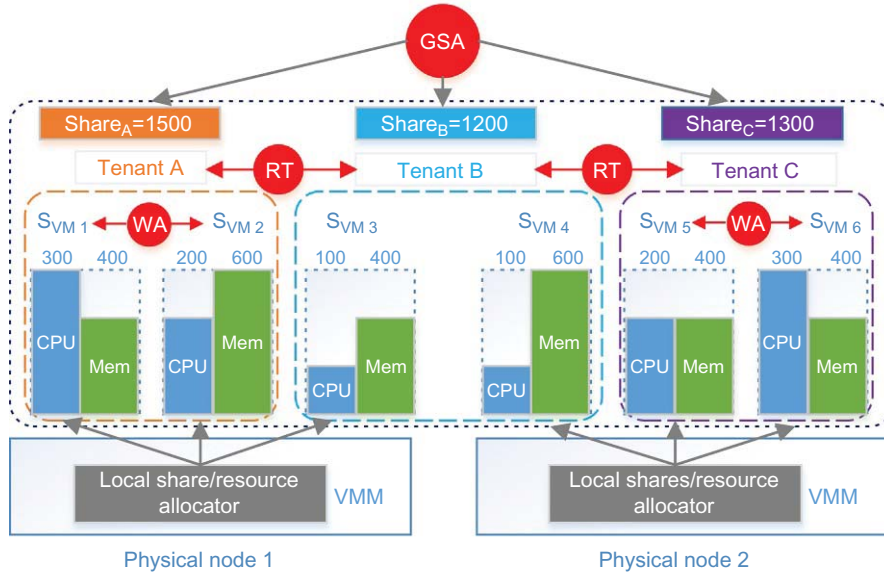


FIG. 2

Hierarchical resource allocation based on resource trading and weight adjustment.

losing underutilized resource. Moreover, tenants can dynamically adjust share allocation of VMs based on actual demands. For example, in Fig. 2, tenant A may deprive 200 memory shares from VM2 and reallocate them to VM1. In this paper, we propose RT between different tenants and dynamic weight adjustment (WA) among multiple VMs belonging to the same tenant. Fig. 2 shows the hierarchical resource allocation based on these two mechanisms. The global share allocator first reserves capacity in bulk based on tenants' aggregate resource demands, then allocates shares to tenants according to their payment. The local share/resource allocator in each node is responsible for RT between tenants, and WA among multiple VMs belonging to the same tenant.

7.5.2 MULTIRESOURCE FAIR SHARING ISSUES

In the following, we demonstrate the deficiency of WMMF and DRF for multiresource sharing in clouds.

Example 1: Assume there are three tenants, each of which has one VM. All VMs share a resource pool consisting of total 20 GHz CPU and 10 GB RAM. Each VM has initial shares for different types of resource when it is created. For example, VM1 initially has CPU and RAM shares of 500 each, simply denoted by a vector $\langle 500, 500 \rangle$. The VMs may have dynamic resource demands. At a time, VM1 runs jobs with demands of 6 GHz CPU and 3 GB RAM, simply denoted by a vector $\langle 6 \text{ GHz}, 3 \text{ GB} \rangle$. The VMs' initial shares and demand vectors are illustrated in Table 3. We examine whether the T-shirt model, WMMF, and DRF can achieve resource efficiency and economic fairness.

With the T-shirt model, we allocate the total resources to tenants in proportion to their share values of CPU and memory separately. The T-shirt model guarantees that each tenant precisely receives the resource shares that the tenant pays for. However, it wastes scarce resource because it may overallocate resource to VMs that has high shares but low demand, even as other VMs have unsatisfied demand. As shown in Table 3, VM2 wastes 1.5 GB RAM and VM3 wastes 2 GHz CPU.

We now apply the WMMF algorithm on each resource type. As shown in Table 3, VM1, VM2, and VM3 initially owns 25%, 25%, 50% of total resource shares, respectively. However, VM1 is allocated with 30% of total resources, with 5% "stolen" from other VMs. Ironically, VM2 contributes 1.5 GB RAM and VM3 contributes 2 GHz CPU to other tenants. However, they do not benefit more than VM1 from resource sharing because CPU and memory resources are allocated separately. In this case, if VM1 deliberately provisions fewer resources than its demand and always counts on others' contributions,

Table 3 Comparison of Resource Allocation Policies Between T-shirt, WMMF, and DRF

VMs	VM1	VM2	VM3	Total
Initial shares	$\langle 500, 500 \rangle$	$\langle 500, 500 \rangle$	$\langle 1000, 1000 \rangle$	$\langle 2000, 2000 \rangle$
Demands	$\langle 6 \text{ GHz}, 3 \text{ GB} \rangle$	$\langle 8 \text{ GHz}, 1 \text{ GB} \rangle$	$\langle 8 \text{ GHz}, 8 \text{ GB} \rangle$	$\langle 22 \text{ GHz}, 12 \text{ GB} \rangle$
T-shirt allocation	$\langle 5 \text{ GHz}, 2.5 \text{ GB} \rangle$	$\langle 5 \text{ GHz}, 2.5 \text{ GB} \rangle$	$\langle 10 \text{ GHz}, 5 \text{ GB} \rangle$	$\langle 18 \text{ GHz}, 8.5 \text{ GB} \rangle$
WMMF allocation	$\langle 6 \text{ GHz}, 3 \text{ GB} \rangle$	$\langle 6 \text{ GHz}, 1 \text{ GB} \rangle$	$\langle 8 \text{ GHz}, 6 \text{ GB} \rangle$	$\langle 20 \text{ GHz}, 10 \text{ GB} \rangle$
WDRF dominant share	$6/20 = 3/10$	$8/20 \text{ CPU}$	$8/(10 * 2) \text{ RAM}$	100%
WDRF allocation	$\langle 6 \text{ GHz}, 3 \text{ GB} \rangle$	$\langle 7 \text{ GHz}, 1 \text{ GB} \rangle$	$\langle 7 \text{ GHz}, 6 \text{ GB} \rangle$	$\langle 20 \text{ GHz}, 10 \text{ GB} \rangle$

then VM1 becomes a free rider. Although WMMF can guarantee resource efficiency, it cannot fully preserve tenant resource shares, and it eventually results in economic unfairness.

We also apply weighted DRF (WDRF) [7] to this example. Both CPU and RAM shares of VM1, VM2, and VM3 correspond to a ratio of 1:1:2. VM1's dominant share can be CPU or memory, both equal to 6/20. VM2's dominant share is CPU share as $\max(8/20, 1/10) = 8/20$. For VM3, its unweighted dominant share is memory share 8/10. Its weight is twice that of VM1 and VM2, so its weighted dominant share is $8/(10 \times 2) = 8/20$. Thus, the ascending order of three VM's dominant shares is $VM1 < VM2 = VM3$. According to WDRF, VM1's demand is first satisfied, then the remanding resources are allocated to VM2 and VM3 based on max-min fairness. We find that VM1 is again a free rider.

In summary, the T-shirt model is not resource efficient, and WMMF and DRF are not economically fair for multiresource allocation. Intuitively, in a cooperative environment, the more one contributes, the more one should gain. Otherwise, tenants would lose their sharing incentives. This is especially important for resource sharing among multiple tenants in pay-as-you-use clouds. Thereby, a new mechanism is needed to reinforce the fairness of multiresource sharing in IaaS clouds.

7.5.3 RECIPROCAL RESOURCE FAIRNESS

In the following, we consider the fair sharing model in a shared system with p types of resource and m tenants. The total system capacity bought by the m tenants is denoted by a vector Ω (ie, $\langle \Omega_{\text{CPU}}, \Omega_{\text{RAM}} \rangle$), denoting the amount of CPU and memory resource, respectively. Each tenant i may have n VMs. Each VM j is initially allocated with a share vector $s(j)$ that reflects its priority relative to other VMs. The amount of resource shares required by VM j is characterized by a demand vector $d(j)$. Correspondingly, the resource share lent to other tenants becomes $s(j) - d(j)$, and we call it a contribution vector $c(j)$. At last, let $s'(j)$ denote the current share vector when resources are reallocated. For simplicity, we assume that resource allocation is oblivious, meaning that the current allocation is not affected by previous allocations. Thus, a VM's priority is always determined by its initial share vector $s(j)$ in each time of resource allocation.

Intertenant resource trading

For multiresource allocation, it is hard to guarantee that the demands of all resource types are nicely satisfied without waste. For example, a tenant's aggregate CPU demand may be less than the initial CPU share, but memory demand exceeds the current allocation. In this case, the tenant may expect to trade CPU resources with other tenants' memory resources. Thus, the question is how to trade resources of different types among tenants while guaranteeing economic fairness. RRF embraces an IRT mechanism with the core idea that tenants' gains from other tenants should be proportional to their contribution. The only basis for underutilized resource allocation is the tenant's contribution, rather than the initial resource share or unsatisfied demand. As shown in Fig. 3, the memory resource contributed by tenant A is twice more than that of tenant B, and tenant A should receive twice more the unused CPU resource (contributed by tenant C) than tenant B at first. Then, we need to check whether the CPU resource of tenant A is overprovisioned. If so, the unused portion should be redistributed to other tenants. This process should be *iteratively* performed by all tenants because each round of resource distribution may affect other tenants' allocations. While this naive approach works, it can cause unacceptable computation overhead. We further propose a work backward strategy to speed up the unused resource distribution.

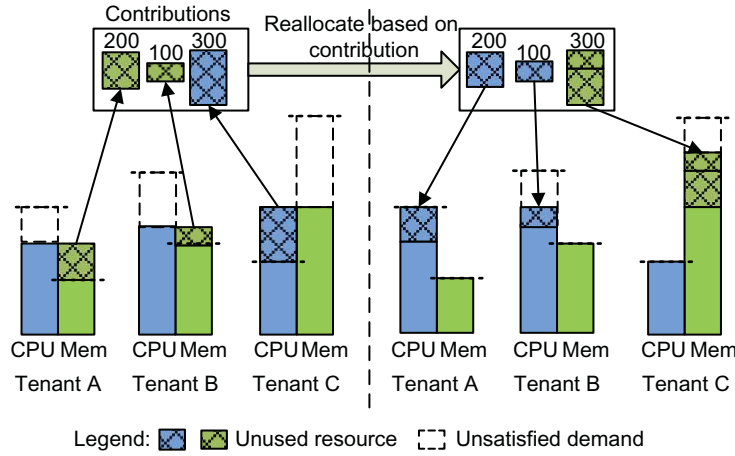


FIG. 3

Sketch of intertenant resource trading.

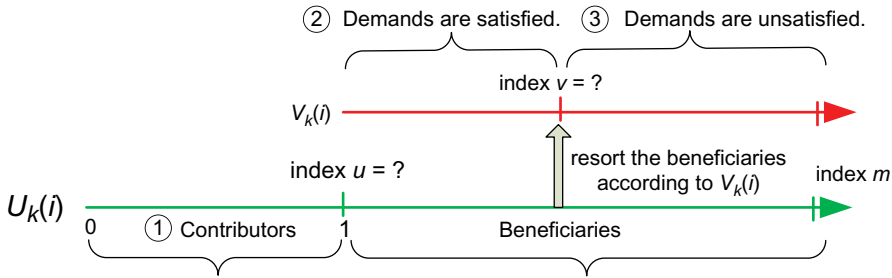


FIG. 4

Sketch of the IRT algorithm.

For each type of resource, we divide the tenants into three categories: contributors, beneficiaries whose demands are satisfied, and beneficiaries whose demands are unsatisfied, as shown in Fig. 4. Tenants in the first two categories are *directly* allocated with their demands exactly, and tenants in the third category are allocated with their initial shares, plus a portion of contributions from the first category. However, a challenging problem is how to divide the tenants into these three categories efficiently. Algorithm 2 describes the sorting process by using some heuristics.

Let vectors $D(i)$, $S(i)$, $C(i)$, and $S'(i)$ denote the total demand, initial share, contribution, and current share of the tenant i , respectively. Correspondingly, let $D_k(i)$, $S_k(i)$, $C_k(i)$, and $S'_k(i)$ denote the total demand, initial share, contribution, and current share of resource type k , respectively. We consider a scenario where m tenants share a resource pool with capacity Ω , with resource contentions (ie, $\sum_{i=1}^m D_i(i) \geq \Omega$). Our algorithm first divide the total capacity on the basis of each tenant's initial share, and then caps each tenant's allocation at the total demand. Actually, each tenant will receive an initial total share S_{piq} , and then the total contribution becomes $C(i) = S(i) - D(i)$ (if $S(i) > D(i)$). For resource type k ($1 \leq k \leq p$), the unused resource $C_k(i)$ is redistributed to other unsatisfied tenants in the ratio of their total contributions.

Algorithm 2 shows the pseudocode for IRT. We first calculate each tenant's total contribution $\Lambda(i)$ (lines 6–8). To reduce the complexity of resource allocation, for each resource type k , we define the normalized demand of tenant i as $U_k(i) = D_k(i) / S_k(i)$, and reindex the tenants so that the $U_k(i)$ are in ascending order, as shown in Fig. 4. Then, we can easily find the index u so that $U_k(u) < 1$ and $U_k(u+1) \geq 1$. The tenants with index $[1, \dots, u]$ are contributors and the remaining are beneficiaries. For tenants with index $[u+1, \dots, m]$ ($U_k(i) \geq 1$), we define the ratio of unsatisfied demand of resource type k to her total contribution as $V_k(i) = (D_k(i) - S_k(i)) / \sum_{k=1}^p C_k(i)$ (lines 12–13), and reindex these tenants according to the ascending order of $V_k(i)$, as shown in Fig. 4. Thus, tenants with index $[1, \dots, u]$ are ordered by $U_k(i)$ while tenants with index $[u+1, \dots, m]$ are ordered by $V_k(i)$. The demand of tenants with the largest index will be satisfied at last. We need to find a pair of successive indexes $v, v+1$ ($v \geq u$), so that the share allocations of tenants with index $[1, \dots, v]$ are capped at their demands, and the remaining contribution $\sum_{i=v+1}^m \Psi_k(i) = \Omega_k - \sum_{i=1}^m D_k(i) - \sum_{i=v+1}^m S_k(i)$ is distributed to tenants with index $[v+1, \dots, m]$ in proportion to their total contributions. Some heuristics can be employed to speed up the index searching. First, searching should start from index $u+1$ because tenants with index $[1, \dots, u]$ are contributors. Second, we can use a binary search strategy to find two successive indices $v, v+1$ that cannot be satisfied. Namely, the following inequality (1) and (2) must be satisfied:

$$S_k(v) + \frac{\sum_{i=v}^m \Psi_k(i) \times \Lambda(v)}{\sum_{i=v}^m \Lambda(i)} \geq D_k(v) \quad (1)$$

$$S_k(v+1) + \frac{\sum_{i=v+1}^m \Psi_k(i) \times \Lambda(v+1)}{\sum_{i=v+1}^m \Lambda(i)} < D_k(v+1) \quad (2)$$

where $\sum_{i=v}^m \Psi_k(i)$ and $\sum_{i=v+1}^m \Psi_k(i+1)$ represent the remaining contributions that will be redistributed to tenants with unsatisfied demands. Once the index v is determined, we can calculate the remaining contribution. The tenants with index $[1, \dots, v]$ receive shares capped at their demands (lines 16–17), and the tenants with index $[v+1, \dots, m]$ receive their initial shares plus the remaining resource in proportion to their contributions (lines 19–20).

ALGORITHM 2 INTERTENANT RESOURCE TRADING (IRT)

Input: $D = \{D(1), \dots, D(m)\}, S = \{S(1), \dots, S(m)\}, \Omega$

Output: $S' = \{S'(1), \dots, S'(m)\}$

Variables: $[i, C(i), \Lambda(i), U(i), V(i), \Psi(i)] \leftarrow 0$

- 1: **for** resource-type $k = 1$ **to** p **do**
- 2: **for** Tenant $i = 1$ **to** m **do**
- 3: /*Allocate each tenant (i) her initial share $S(i)$ */
- 4: $S'_k(i) \leftarrow S_k(i)$
- 5: $U_k(i) \leftarrow D_k(i) / S_k(i)$
- 6: **if** $S_k(i) \geq D_k(i)$ **then**
- 7: $C_k(i) \leftarrow S_k(i) - D_k(i)$

```

/*Calculate tenant(i)'s total contribution  $\Lambda(i)$  on all type of resource */
8:  $\Lambda(i) \leftarrow \Lambda(i) + C_k(i)$ 
9: for resource-type  $k = 1$  to  $p$  do
10: Sort  $U_k(i)$  in ascending order;
11: Find the index  $u$  so that  $U_k(k) < 1 \leq U_k(u+1)$ 
12: for Tenant  $i = u+1$  to  $m$  do
13:  $V_k(i) \leftarrow (D_k(i) - S_k(i)) / \Lambda(i)$ 
14: Sort  $V_k(i)$  in ascending order;
15: Find the index  $v$  using binary search algorithm so that Equation (1) and (2) are satisfied;
16: for Tenant  $i = 1$  to  $v$  do
17:  $S'_k(i) \leftarrow D_k(i)$  /*share is capped by demand*/
/*Calculate the remaining contributions for re-allocation*/
18:  $\sum_{i=v+1}^m \Psi_k(i) \leftarrow \Omega_k - \sum_{i=1}^v D_k(i) - \sum_{i=v+1}^m S_k(i)$ 
19: for Tenant  $i = v+1$  to  $m$  do
20:  $S'_k(i) \leftarrow S_k(i) + \frac{\sum_{i=v+1}^m \Psi_k(i) \times \Lambda(v+1)}{\sum_{i=v+1}^m \Lambda(i)}$ 

```

Intratenant weight adjustment

A tenant usually needs more than one VM to host her applications. Workloads in different VMs may have dynamic and heterogeneous resource requirements. Thus, dynamic resource flows among VMs belonging to the same tenant can prevent loss of tenant's asset. In F2C, we use IWA to adjust the resource among VMs belonging to the same tenant. We allocate *share* (or weight) for each VM using a policy similar to WMMF. For each type of resource, we first reset each VM's current weight to its initial share. However, if the allocation made to a VM is more than its demand, its allocation should be capped at its real demand, and the unused share should be reallocated to its sibling VMs with unsatisfied demands. In contrast to WMMF that reallocates the unused resource in proportion to VMs' share values, we reallocate the excessive resource share to the VMs in the ratio of their unsatisfied demands.

Note that once a VM's resource share is determined, the resource allocation made to the VM is simply determined by the function $share \xrightarrow{f_2} resource$.

ALGORITHM 3 INTERTENANT WEIGHT ADJUSTMENT (IWA)

Input: $d = \{d(1), \dots, d(n)\}, s = \{s(1), \dots, s(n)\}, S$

Output: $s' = \{s'(1), \dots, s'(n)\}$

Variables: $[j, \Gamma, \Phi] \leftarrow 0$

/* Allocate initial share $s(j)$ to each VM(j) */

```

1:  $\Phi \leftarrow S - \sum_{j=1}^n s(j)$  /*Calculate the difference of initial total share and new allocated capacity */
2: for VM  $j = 1$  to  $n$  do
3:   if  $d(j) \geq s(j)$  then
4:      $\Gamma \leftarrow \Gamma + (d(j) - s(j))$  /*total unsatisfied demand*/
5:   else
6:      $\Phi \leftarrow \Phi + (s(j) - d(j))$  /*total remaining capacity */
/* distribute remaining capacity to VMs with unsatisfied demand */
7: for VM  $j = 1$  to  $n$  do
8:   if  $d(j) \geq s(j)$  then
9:      $s'(j) \leftarrow s(j) + \frac{d(j) - s(j)}{\Gamma} \times \Phi$ 
10:  else
11:     $s'(j) \leftarrow d(j)$ 

```

Algorithm 3 shows the pseudocode for IWA. A tenant with n VMs is allocated with total resource share S . Note that S is a global allocation vector, which corresponds to the output of **Algorithm 2**. Thus, **Algorithm 3** is performed by accompanying **Algorithm 2**. For each tenant, we first calculate the total unsatisfied demand and total remaining capacity for reallocation, respectively (lines 2–6), then distribute the remaining capacity to unsatisfied VMs in a ratio of their unsatisfied demands (lines 7–11). As VM provisioning is constrained to physical hosts' capacity, it is desirable to adjust weights of VMs on the same physical node, rather than across multiple nodes. In practice, we execute the IWA algorithm only on each single node.

7.5.4 EXPERIMENTAL EVALUATION

RRF is implemented on top of Xen 4.1. The prototype F2C is deployed in a cluster with 10 nodes. The following workloads with diversifying and variable resource requirements are used to evaluate resource allocation fairness and application performance.

TPC-C: We use a public benchmark DBT-2 as clients to simulate a complete computing environment where a number of users execute transactions against a database. The clients and MySQL database server run in two VMs separately. We assume these two VM belong to the same tenant. We configure 600 terminal threads and 300 database connections. We evaluate its application performance by throughput (transactions per minute).

RUBBoS (Rice University Bulletin Board System) [38]: it is a typical multitier web benchmark. RUBBoS simulates an online news forum like <https://slashdot.org>. We deploy the benchmark in a three-tier architecture using Apache 2.2, Tomcat 7.0, and MySQL 5.5.

Kernel-build: We compile Linux 2.6.31 kernel in a single VM. It generates a moderate and balanced load of CPU and memory.

Hadoop: We use Hadoop WordCount microbenchmark to setup a virtual cluster consisting of 10 worker VMs and 1 master VM. The 10 worker VMs are evenly distributed in the 10 physical machines and co-run with other three workloads.

On our testbed, we consider multiple tenants sharing the cluster. Each tenant only runs one kind of the above workloads. All workloads are running in VMs whose initial share values are set according to the workloads' average demands. In addition, we can configure the initial share of tenant (i) based on a provisioning coefficient:

$$\alpha = S(i) / \overline{D(I)}$$

which reflects the ratio of initial resource share to the average demand. We continuously launch the tenants' applications to the cluster one by one until there is no room to accommodate any more applications.

We evaluate F2C by comparing the following alternative approaches for IaaS clouds:

- *T-shirt (static)*: Workloads are running in VMs with static resource provision. It is the current resource model adopted by most IaaS clouds [1].
- *WMMF*: WMMF [35] is used to allocate CPU and memory resources to each VM separately.
- *DRF*: DRF [7] is used to allocate multiple resources to each VM.
- *IWA*: We conduct only weight adjustment for VMs belonging to the same tenant, without considering IRT. This is to assess the individual impact of IRT.
- *RRF (IRT + IWA)*: We conduct hierarchical resource allocation using both IRT and IWA.

Results on fairness

In general, in a shared computing system with resource contention, tenants want to receive more resources or at least the same amount of resource than they buy. We call it *fair* if a tenant can achieve this goal (ie, sharing benefit). In contrast, it is also possible that the total resources tenants receives is fewer than that without sharing, which we call *unfair* (ie, sharing loss). Thus, we define the economic fairness degree $\beta(i)$ for tenant i in a time window T as follows:

$$\beta(i) = \sum_{i=1}^T S'_i(i) / (T \times S(i))$$

It represents the ratio of average resource shares received to the tenant's initial shares in the time window T ; more exactly speaking, it denotes the ratio of resource value to tenants payments. $\beta(i) = 1$ implies absolute economic fairness. $\beta(i) > 1$ implies the tenant benefits from resource sharing, while $\beta(i) < 1$ implies the tenant loses her asset.

We evaluate the economic fairness of different schemes for the four workloads in a long period of time. All VMs' capacities are provisioned based on their average demands (ie, $\alpha = 1$). Fig. 5 shows the comparison of economic fairness of different resource allocation schemes. Each bar shows the average result of the same workload run by multiple tenants. Overall, RRF achieves much better economic fairness than other approaches. Specifically, RRF leads to smaller difference of β between different applications, indicating 95% economic fairness (geometric mean) for multiresource sharing among multitenants.

We make the following observations:

First, the T-shirt (static) model achieves 100% economic fairness, as VMs share nothing with each other. However, it results in the worst application performance for all workloads, as shown in Fig. 6.

Second, both WMMF and DRF show significant differences on β for different workloads. As all WMMF-based algorithms always try to satisfy the demand of smallest applications first, both kernel-build and TPC-C gain more resources than their initial shares. This effect is more significant for DRF if the application shows a tiny skewness of multiresource demands (such as kernel-build). DRF always completely satisfies the demand of these applications first. Thus, DRF and WMMF are susceptible to application's

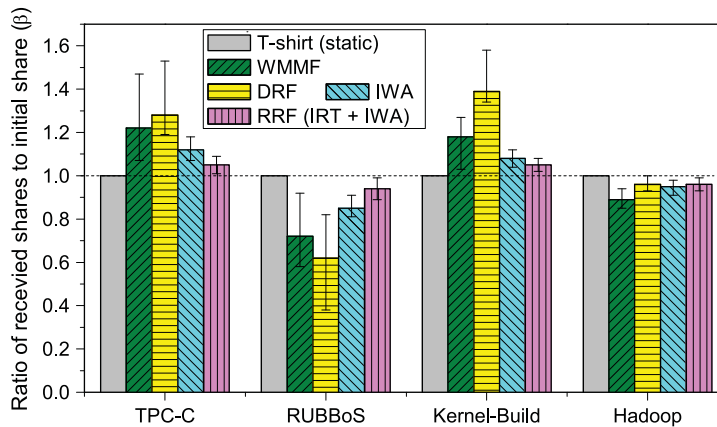


FIG. 5

Comparison of the economic fairness of several resource allocation schemes.

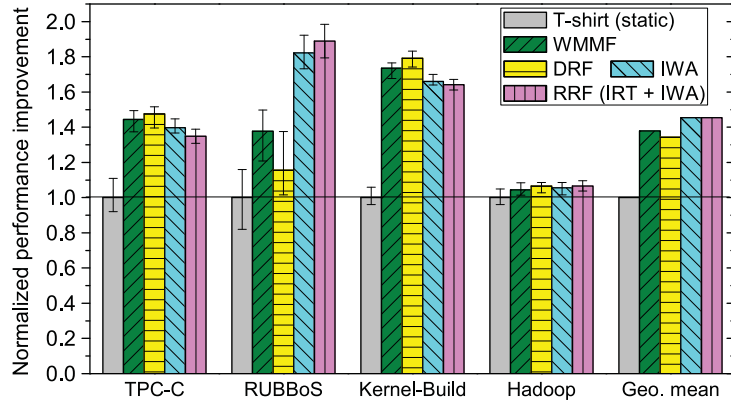


FIG. 6

Comparison of the application performance improvement of several resource allocation schemes.

load patterns. Applications with large skewness of multiresource demands usually lose their asset. Even for a single-resource type, a large deviation of resource demand also lead to distinct economic unfairness.

Third, workloads with different resource demand patterns show different behaviors in resource sharing. RUBBoS has a cyclical workload pattern, its resource demand shows the largest temporal fluctuations. When the load is below its average demand, it contributes resource to other VMs and thus loses its asset. However, when it shows a large value of D_i , $\pi_i \{S_{piq}$, its demand always cannot be fully satisfied if there exists resource contention.

Thus, RUBBoS has smaller value of β than other applications. For TPC-C and kernel-build, although they show comparable demand deviation and skewness with RUBBoS, they have much more opportunities to benefit from RUBBoS because their absolute demands are much smaller. For Hadoop, although it requires a large amount of resources, it demonstrates only a slight deviation of resource demands, and there are few opportunities to contribute resources to other VMs (except in its *reduce* stage).

Fourth, IWA allows tenants to properly distribute VMs' spare resources to their sibling VMs in proportion to their unsatisfied demands. It guarantees that tenants can effectively utilize their own resource. IRT can further preserve tenants' assets as each tenant tries to maximize the value of spare resource. In addition, RRF is immune to free riding.

We also studied the impact of different α values. For space limitation, we omit the figures and briefly discuss the results. When we decrease the provisioning coefficient α (ie, reduce the resource provisioned for applications), the β of all applications approaches to one. In contrast, a larger α leads to a decrease of β . That means, applications tend to preserve their resources when there exist intensive resource contention. Nevertheless, a larger value of α implies better application performance.

Improvement of application performance

Fig. 6 shows the normalized application performance for different resource allocation schemes. All schemes provision resources according to the applications' average demand ($\alpha=1$). In the T-shirt model, all applications show the worst performance and we refer it as a baseline. In other models, all applications show performance improvement due to resource sharing. For RUBBoS, RRF leads to much more improved application performance than other schemes. This is because RRF provides

two mechanisms (IRT + IWA) to preserve tenants' assets, and thus RRF allow RUBBoS receive more resources than other schemes, as shown in Fig. 5. For other workloads, RRF is also comparable to the other resource-sharing schemes. In summary, RRF achieves 45% performance improvement for all workloads on average (geometric mean). DRF achieves the best performance for kernel-build and TPC-C, but achieves very bad performance for RUBBoS. It shows the largest performance differentiation for different workloads. DRF always tends to satisfy the demand of the application with the smallest dominant share, and thus applications that have resource demands of small sizes or small skewness always benefit more from resource sharing. In contrast, the performance of Hadoop shows slight variations between different allocation schemes due to its rather stable resource demands.

7.6 RELATED WORK ON RESOURCE MANAGEMENT

In this section, we review the related work of resource management from the aspects of performance, fairness, energy, and power cost, as well as monetary cost. For more related work on resource management, we refer readers to three comprehensive surveys [39–41].

7.6.1 RESOURCE UTILIZATION OPTIMIZATION

For resource management, different resource allocation strategies can lead to significantly varied resource utilization [42]. Various resource allocation approaches have therefore been proposed for different workloads and systems/platforms [43–46]. For MapReduce workloads, Tang et al. [45] observed that different job submission order can have a significant impact on the resource utilization, and therefore improved the resource utilization for MapReduce cluster by reordering the job submission order of arriving MapReduce jobs. Grandl et al. [44] proposed a dynamic multiresource packing system called Tetris that improves the resource utilization by packing tasks to machines based on their requirements along multiple resources. Moreover, resource sharing is another approach to improve the resource utilization in a multitenant system by allowing overloaded users to possess unused resources from underloaded users [47]. In addition, task/VM migration and consolidation [48–50], widely used in cloud computing, is also an efficient method to improve resource utilization of a single machine.

7.6.2 POWER AND ENERGY COST SAVING OPTIMIZATION

Power and energy are a big concern in current data centers and supercomputers, which consists of hundreds of servers. Efficient resource management is nontrivial for power and energy cost saving. There are a number of techniques proposed to alleviate it. One intuitive approach is shutting down some low-utilized servers [51,52]. Moreover, VM migration and consolidation is an effective approach to reduce the number of running machines and, in turn, save the power and energy cost [53,54]. Finally, we refer readers to a survey [55] for more solutions.

7.6.3 MONETARY COST OPTIMIZATION

Monetary cost optimization has become a hot topic in recent years, especially for cloud computing. A lot of job scheduling and resource provisioning algorithms have been proposed by leveraging

market-based techniques [56], rule-based techniques [57], and model-based approaches [58]. Many relevant cost optimization approaches can be found in databases [59], Internet [60], distributed systems [61], grid [17], and cloud [56].

7.6.4 FAIRNESS OPTIMIZATION

Fairness is an important issue in a multiuser computing environment. There are various kinds of fair policies in the traditional HPC and grid computing, including round-robin [62], proportional resource sharing [63], weighted fair queuing [64], and max-min fairness [65]. In comparison, max-min fairness is the most popular and widely used policy in many existing parallel and distributed systems, such as Hadoop [66], YARN [15], Mesos [16], Choosy [67], and Quincy [68]. Hadoop [66] partitions resources into slots and allocates them fairly across pools and jobs. In contrast, YARN [15] divides resources into containers (ie, a set of various resources like memory and CPU) and tries to guarantee fairness between queues. Mesos [16] enables multiple diverse computing frameworks such as Hadoop and Spark sharing a single system. It proposes a distributed two-level scheduling mechanism called resource offers, which decides how many resources to offer. Each framework decides which resources to accept or which computation to run on them. Choosy [67] extends the max-min fairness by considering placement constraints. Quincy [68] is a fair scheduler for Dryad that achieves a fair scheduling of multiple jobs by formulating it as a min-cost flow problem. In addition to the single-resource fairness, there are some work focusing on multiresource fairness, including DRF [7] and its extensions [69–72].

7.7 OPEN PROBLEMS

Despite many recent efforts on resource management, there are a number of open problems remained to be explored in future. We briefly elaborate on them from the following aspects.

7.7.1 SLA GUARANTEE FOR APPLICATIONS

In practice, users' applications are often with different performance (eg, latency, throughput) and resource requirements. For example, latency-sensitive applications in memcached require a high response time, whereas batch jobs in Hadoop are often require high throughput [73]. When these applications run together in a shared computing system, it becomes challenging work to provide the quality of services for each application. Despite the many research efforts that have been made [74,75], there is a lack of systematic approach that takes into account different types of resources and application requirements integrally. Most of them either focus on a specific resource type (eg, network flow) [73,74] or a kind of application [76]. Existing works do not address how to systematically ensure SLA guarantee for different applications.

7.7.2 VARIOUS COMPUTATION MODELS AND SYSTEMS

As listed in Section 7.3, there are a number of computation models, as well as computing systems proposed for different applications in recent years. From a user's perspective, it becomes a headache and a time-consuming problem for the user to choose and learn those computation models and

corresponding computing systems. Designing a general computing and resource management system like an operation system becomes a challenging issue.

7.7.3 EXPLOITING EMERGING HARDWARE

Emerging hardware is available at different layers. For example, in the storage layer, we now have solid-state disk and nonvolatile RAM, which are much faster than Hard Disk (HD). In the computation layer, there are a set of accelerators such as GPU, AMD Accelerated Processing Unit (APU), and Field Programmable Gate Array (FPGA). Moreover, in the network layer, remote direct memory access is an efficient hardware tool for speeding network transfer. For a computing system, it is important to adopt this emerging hardware to improve the performance of applications. Currently, the study on this aspect is still at early stage. More research efforts are required to efficiently utilize this emerging hardware at different layers for existing computing systems.

7.8 SUMMARY

In this chapter, we have discussed the importance of resource management for Big Data processing and surveyed a number of existing representative large-scale data processing systems. One of the classic issues for resource management is fairness. The chapter reviewed the memory less fair resource allocation policies for existing systems and showed their unsuitability for cloud computing by presenting three problems. A new LTRF policy was then proposed to address these problems, and we formally and experimentally validate the merits of the proposed policy. This chapter next focused on the resource management for VMs on the cloud, considering VM migration and consolidation in the cloud environment. Finally, there are many open problems that need more research efforts in this field.

REFERENCES

- [1] Gmach D, Rolia J, Cherkasova L. Selling t-shirts and time shares in the cloud. In: Proceedings of the 2012 12th IEEE/ACM international symposium on cluster, cloud and grid computing (CCGRID 2012), CCGRID'12. Washington, DC: IEEE Computer Society; 2012. p. 539–46.
- [2] AutoScaling. <http://aws.amazon.com/autoscaling>.
- [3] Mao M, Humphrey M. A performance study on the VM startup time in the cloud. In: IEEE 5th international conference on cloud computing (CLOUD). Washington, DC: IEEE Computer Society; 2012. p. 423–30.
- [4] Nguyen H, Shen Z, Gu X, Subbiah S, Wilkes J. Agile: elastic distributed resource scaling for infrastructure-as-a-service. In: USENIX ICAC; San Jose, CA: USENIX Association; 2013. pp. 69–82.
- [5] Amazon. <http://aws.amazon.com/solutions/case-studies/>.
- [6] Google. <https://cloud.google.com/customers/>.
- [7] Ghodsi A, Zaharia M, Hindman B, Konwinski A, Shenker S, Stoica I. Dominant resource fairness: fair allocation of multiple resource types. In: Proceedings of the 8th USENIX conference on networked systems design and implementation, NSDI'11. Berkeley, CA: USENIX Association; 2011. p. 24.
- [8] Barham P, Dragovic B, Fraser K, Hand S, Harris T, Ho A, et al. Xen and the art of virtualization. In: Proceedings of the nineteenth ACM symposium on operating systems principles, SOSP'03. New York, NY: ACM; 2003. p. 164–77.

- [9] Cherkasova L, Gupta D, Vahdat A. Comparison of the three CPU schedulers in Xen. *ACM SIGMETRICS Perform Eval Rev* 2007;35(2):42–51.
- [10] Gulati A, Merchant A, Varman PJ. mclock: handling throughput variability for hypervisor io scheduling. In: *Proceedings of the 9th USENIX conference on operating systems design and implementation, OSDI'10*. Berkeley, CA: USENIX Association; 2010. p. 1–7.
- [11] Waldspurger CA. Memory resource management in VMware ESX server. In: *Proceedings of the 5th symposium on operating systems design and implementation — copyright restrictions prevent ACM from being able to make the PDFs for this conference available for downloading, OSDI'02*; New York, NY: ACM; 2002. p. 181–94.
- [12] Dean J, Ghemawat S. MapReduce: simplified data processing on large clusters. *Commun ACM* 2008;51(1):107–13.
- [13] Zaharia M, Chowdhury M, Franklin MJ, Shenker S, Stoica I. Spark: cluster computing with working sets. In: *Proceedings of the 2nd USENIX conference on hot topics in cloud computing, HotCloud'10*. Berkeley, CA: USENIX Association; 2010. p. 10.
- [14] Yoo AB, Jette MA, Grondona M. SLURM: simple Linux utility for resource management. In: *Job scheduling strategies for parallel processing*. Seattle, WA: Springer; 2003. p. 44–60.
- [15] Vavilapalli VK, Murthy AC, Douglas C, Agarwal S, Konar M, Evans R, et al. Apache Hadoop YARN: yet another resource negotiator. In: *Proceedings of the 4th annual symposium on cloud computing, SOCC'13*. New York, NY: ACM; 2013. p. 15–6.
- [16] Hindman B, Konwinski A, Zaharia M, Ghodsi A, Joseph AD, Katz R, et al. Mesos: a platform for fine-grained resource sharing in the data center. In: *Proceedings of the 8th USENIX conference on networked systems design and implementation, NSDI'11*. Berkeley, CA: USENIX Association; 2011. p. 22.
- [17] Storage resource management. https://en.wikipedia.org/wiki/Storage_Resource_Management.
- [18] Software-defined networking. <https://en.wikipedia.org/wiki/Software-defined-networking>.
- [19] Malewicz G, Austern MH, Bik AJ, Dehnert JC, Horn I, Leiser N, et al. Pregel: a system for large-scale graph processing. In: *Proceedings of the 2010 ACM SIGMOD international conference on management of data, SIGMOD'10*. New York, NY: ACM; 2010. p. 135–46.
- [20] Hadoop. <http://hadoop.apache.org/>.
- [21] Isard M, Budiu M, Yu Y, Birrell A, Fetterly D. Dryad: distributed data-parallel programs from sequential building blocks. In: *Proceedings of the 2nd ACM SIGOPS/EuroSys European conference on computer systems 2007, EuroSys'07*. New York, NY: ACM; 2007. p. 59–72.
- [22] Giraph. <http://giraph.apache.org/>.
- [23] Storm. <http://storm-project.net/>.
- [24] Tez. <https://tez.apache.org/>.
- [25] Zaharia M, Chowdhury M, Das T, Dave A, Ma J, McCauley M, et al. Resilient distributed datasets: a fault-tolerant abstraction for in-memory cluster computing. In: *Proceedings of the 9th USENIX conference on networked systems design and implementation, NSDI'12*. Berkeley, CA: USENIX Association; 2012. p. 2.
- [26] Olston C, Reed B, Srivastava U, Kumar R, Tomkins A. Pig latin: a not-so-foreign language for data processing. In: *Proceedings of the 2008 ACM SIGMOD international conference on management of data*. New York, NY: ACM; 2008. p. 1099–110.
- [27] Thusoo A, Sarma JS, Jain N, Shao Z, Chakka P, Anthony S, et al. Hive: a warehousing solution over a map-reduce framework. *Proc VLDB Endowment* 2009;2(2):1626–9.
- [28] George L. HBase: the definitive guide. Newton, MA: O'Reilly Media; 2011.
- [29] Xin RS, Rosen J, Zaharia M, Franklin MJ, Shenker S, Stoica I. Shark: SQL and rich analytics at scale. In: *Proceedings of the 2013 ACM SIGMOD international conference on management of data*. New York, NY: ACM; 2013. p. 13–24.

- [30] Armbrust M, Xin RS, Lian C, Huai Y, Liu D, Bradley JK, et al. Spark SQL: relational data processing in spark. In: Proceedings of the 2015 ACM SIGMOD international conference on management of data. New York, NY: ACM; 2015. p. 1383–94.
- [31] Gonzalez JE, Xin RS, Dave A, Crankshaw D, Franklin MJ, Stoica I. Graphx: Graph processing in a distributed dataflow framework. In: Proceedings of OSDI; Broomfield, CO: USENIX Association 2014. p. 599–613.
- [32] Zhong J, He B. Medusa: simplified graph processing on gpus. *IEEE Trans Parallel Distrib Syst* 2014;25(6):1543–52.
- [33] Loan agreement. http://en.wikipedia.org/wiki/Loan_agreement.
- [34] Faraz A, Seyong L, Mithuna T, Puma VTN. Purdue MapReduce benchmarks suite. Technical report EECS-2012, October, West Lafayette, IN: School of Electrical and Computer Engineering, Purdue University; 2012.
- [35] Keshav S. An engineering approach to computer networking: atm networks, the internet, and the telephone network. Boston, MA: Addison-Wesley; 1997.
- [36] Liu H, He B. Reciprocal resource fairness: towards cooperative multiple-resource fair sharing in IaaS clouds. In: International conference for high performance computing, networking, storage and analysis, SC14, November; New Orleans, LA: IEEE; 2014. p. 970–81.
- [37] Williams D, Jamjoom H, Liu Y-H, Weatherspoon H. Overdriver: handling memory overload in an oversubscribed cloud. In: Proceedings of the 7th ACM SIG-PLAN/SIGOPS international conference on virtual execution environments, VEE'11. New York, NY: ACM; 2011. p. 205–16.
- [38] RUBBoS. <http://jmob.ow2.org/rubbos.html>.
- [39] Jennings B, Stadler R. Resource management in clouds: survey and research challenges. *J Netw Syst Manag* 2014;23(3):1–53.
- [40] Krauter K, Buyya R, Maheswaran M. A taxonomy and survey of grid resource management systems for distributed computing. *Softw Pract Exper* 2002;32(2):135–64.
- [41] Weingrtner R, Brscher GB, Westphall CB. Cloud resource management: a survey on forecasting and profiling models. *J Netw Comput Appl* 2015;47:99–106.
- [42] Anuradha V, Sumathi D. A survey on resource allocation strategies in cloud computing. In: International conference on IEEE information communication and embedded systems (ICICES) Chennai, TN: IEEE; 2014. p. 1–7.
- [43] Delimitrou C, Kozyrakis C. Quasar: resource-efficient and QoS-aware cluster management. *ACM SIGPLAN Not* 2014;49(4):127–44.
- [44] Grandl R, Ananthanarayanan G, Kandula S, Rao S, Akella A. Multi-resource packing for cluster schedulers. In: Proceedings of the 2014 ACM conference on SIGCOMM, SIGCOMM'14. New York, NY: ACM; 2014. p. 455–66.
- [45] Tang S, Lee B-S, He B. MROrder: flexible job ordering optimization for online MapReduce workloads. In: Wolf F, Mohr B, Mey D, editors. Euro-Par 2013 parallel processing. Lecture notes in computer science, vol. 8097. Berlin, Heidelberg: Springer; 2013. p. 291–304.
- [46] Verma A, Pedrosa L, Korupolu M, Oppenheimer D, Tune E, Wilkes J. Large-scale cluster management at Google with borg. In: Proceedings of the tenth European conference on computer systems. New York, NY: ACM; 2015. p. 18.
- [47] Tang S, Lee B-s, He B, Liu H. Long-term resource fairness: towards economic fairness on pay-as-you-use computing systems. In: Proceedings of the 28th ACM international conference on supercomputing, ICS'14. New York, NY: ACM; 2014. p. 251–60.
- [48] Beloglazov A, Buyya R. Managing overloaded hosts for dynamic consolidation of virtual machines in cloud data centers under quality of service constraints. *IEEE Trans Parallel Distrib Syst* 2013;24(7):1366–79.
- [49] Corradi A, Fanelli M, Foschini L. VM consolidation: a real case based on openstack cloud. *Futur Gener Comput Syst* 2014;32:118–27.
- [50] Hsu C-H, Chen S-C, Lee C-C, Chang H-Y, Lai K-C, Li K-C, et al. Energy-aware task consolidation technique for cloud computing. In: IEEE third international conference on cloud computing technology and science (CloudCom). Washington, DC: IEEE; 2011. p. 115–21.

- [51] Lang W, Patel JM. Energy management for MapReduce clusters. *Proc VLDB Endowment* 2010;3(1-2):129–39.
- [52] Rajamani K, Lefurgy C. On evaluating request-distribution schemes for saving energy in server clusters. In: *IEEE international symposium on performance analysis of systems and software, ISPASS, 2003*. Washington, DC: IEEE; 2003. p. 111–22.
- [53] Clark C, Fraser K, Hand S, Hansen JG, Jul E, Limpach C, et al. Live migration of virtual machines. In: *Proceedings of the 2nd conference on symposium on networked systems design & implementation*. vol. 2. Berkeley, CA: USENIX Association; 2005. p. 273–86.
- [54] Ranganathan P, Leech P, Irwin D, Chase J. Ensemble-level power management for dense blade servers. *ACM SIGARCH Computer Architecture News* 2006;34:66–77.
- [55] Beloglazov A, Buyya R, Lee YC, Zomaya A, et al. A taxonomy and survey of energy-efficient data centers and cloud computing systems. *Adv Comput* 2011;82(2):47–111.
- [56] Fard HM, Prodan R, Fahringer T. A truthful dynamic workflow scheduling mechanism for commercial multicloud environments. *IEEE Trans Parallel Distrib Syst* 2013;24(6):1203–12.
- [57] Malawski M, Juve G, Deelman E, Nabrzyski J. Cost-and deadline-constrained provisioning for scientific workflow ensembles in IaaS clouds. In: *Proceedings of the international conference on high performance computing, networking, storage and analysis*. Los Alamitos, CA: IEEE Computer Society Press; 2012. p. 22.
- [58] Byun E-K, Kee Y-S, Kim J-S, Maeng S. Cost optimized provisioning of elastic resources for application workflows. *Future Generation Computer Systems* 2011;27(8):1011–26.
- [59] Gray J, Graefe G. The five-minute rule ten years later, and other computer storage rules of thumb. *ACM SIGMOD Record*. New York, NY: ACM. 1997;26(4):63–68.
- [60] Ma RT, Chiu DM, Lui J, Misra V, Rubenstein D. Internet economics: the use of Shapley value for ISP settlement. *IEEE/ACM Trans Networking* 2010;18(3):775–87.
- [61] Gray J. Distributed computing economics. *Queue* 2008;6(3):63–8.
- [62] Drozdowski M. *Scheduling for parallel processing*. 1st ed. New York, NY: Springer; 2009.
- [63] Waldspurger CA, Weihl WE. Lottery scheduling: flexible proportional-share resource management. In: *Proceedings of the 1st USENIX conference on operating systems design and implementation, OSDI'94*. Berkeley, CA: USENIX Association; 1994.
- [64] Demers A, Keshav S, Shenker S. Analysis and simulation of a fair queueing algorithm. In: *Symposium proceedings on communications architectures & protocols, SIGCOMM'89*. New York, NY: ACM; 1989. p. 1–12.
- [65] Max-min fairness (wikipedia). http://en.wikipedia.org/wiki/Max-min_fairness.
- [66] White T. *Hadoop: the definitive guide*. 1st ed. Sebastopol, CA: O'Reilly; 2009. June.
- [67] Ghodsi A, Zaharia M, Shenker S, Stoica I. Choosy: max-min fair sharing for datacenter jobs with constraints. In: *Proceedings of the 8th ACM European conference on computer systems, EuroSys'13*; New York, NY: ACM; 2013. p. 365–78.
- [68] Isard M, Prabhakaran V, Currey J, Wieder U, Talwar K, Goldberg A. Quincy: fair scheduling for distributed computing clusters. In: *Proceedings of the ACM SIGOPS 22nd symposium on operating systems principles, SOSP'09*. New York, NY: ACM; 2009. p. 261–76.
- [69] Bhattacharya AA, Culler D, Friedman E, Ghodsi A, Shenker S, Stoica I. Hierarchical scheduling for diverse datacenter workloads. In: *Proceedings of the 4th annual symposium on cloud computing, SOCC'13*. New York, NY: ACM; 2013. p. 14–5.
- [70] Kash I, Procaccia AD, Shah N. No agent left behind: Dynamic fair division of multiple resources. In: *Proceedings of the 2013 international conference on autonomous agents and multi-agent systems, AAMAS'13*. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems; 2013. p. 351–8.
- [71] Parkes DC, Procaccia AD, Shah N. Beyond dominant resource fairness: extensions, limitations, and indivisibilities. In: *Proceedings of the 13th ACM conference on electronic commerce, EC'12*. New York, NY: ACM; 2012. p. 808–25.

- [72] Wang W, Li B, Liang B. Dominant resource fairness in cloud computing systems with heterogeneous servers. In: Proceedings IEEE INFOCOM, April 2014. Toronto, ON: IEEE; 2014. p. 583–91.
- [73] Grosvenor MP, Schwarzkopf M, Gog I, Watson RN, Moore AW, Hand S, et al. Queues dont matter when you can jump them!. In: NSDI'15 proceedings of the 12th USENIX conference on networked systems design and implementation. Berkeley, CA: USENIX Association; 2015. p. 1–14.
- [74] Chowdhury M, Stoica I. Coflow: an application layer abstraction for cluster networking. Proceedings of the 11th ACM Workshop on Hot Topics in Networks. New York, NY: ACM Hotnets; 2012. pp. 31–36. ACM.
- [75] Mace J., Bodik P., Fonseca R., Musuvathi M.. Retro: targeted resource management in multi-tenant distributed systems. In: Proceedings of the 12th USENIX Conference on Networked Systems Design and Implementation (NSDI'15). Berkeley, CA: USENIX Association; 2015. p. 589–603.
- [76] Lim N, Majumdar S, Ashwood-Smith P. A constraint programming-based resource management technique for processing MapReduce jobs with SLAS on clouds. In: 43rd international conference on IEEE parallel processing (ICPP). Minneapolis, MN: IEEE; 2014. p. 411–21.