

OPERA: Opportunistic and Efficient Resource Allocation in Hadoop YARN by Harnessing Idle Resources

Yi Yao*

yyao@ece.neu.edu

Han Gao*

hgao@ece.neu.edu

Jiayin Wang[†]

jane@cs.umb.edu

Ningfang Mi*

ningfang@ece.neu.edu

Bo Sheng[†]

shengbo@cs.umb.edu

*Department of Electrical and Computer Engineering, Northeastern University, 360 Huntington Ave., Boston, MA 02115

[†]Department of Computer Science, University of Massachusetts Boston, 100 Morrissey Boulevard, Boston, MA 02125

Abstract—Efficiently managing resources to improve the throughput in a large-scale cluster has become a crucial problem with the explosion of data processing applications in recent years. Hadoop YARN and Mesos, as two universal resource management platforms, have been widely adopted in the commodity cluster for co-deploying multiple data processing frameworks, such as Hadoop MapReduce and Apache Spark. However, in the existing resource management, a certain amount of resources are *exclusively* allocated to a running task and can only be re-assigned after that task is completed. This exclusive mode unfortunately may underutilize the cluster resources and degrade system performance. To address this issue, we propose a novel opportunistic and efficient resource allocation approach, named **OPERA**, which breaks the barriers among the encapsulated resource containers by leveraging the knowledge of actual runtime resource utilizations to re-assign opportunistic available resources to the pending tasks. We implement and evaluate **OPERA** in Hadoop YARN v2.5. Our experimental results show that **OPERA** significantly reduces the average job execution time and increases the resource (CPU and memory) utilizations.

Keywords: Resource allocation, MapReduce scheduling, Hadoop YARN.

I. INTRODUCTION

Large-scale data processing has become ubiquitous in the era of big data. Many cluster computing frameworks have been developed to simplify distributed data processing on clusters of commodity servers in the past decades. Hadoop MapReduce [1], [2], as one of the prominent frameworks, has been widely adopted in both academia and industry for un-structured data processing [3]. As data sources become more diverse, new frameworks are emerging in recent years and thriving to address different large-scale data processing problems. For example, Apache Spark [4], [5] was introduced to optimize iterative data processing and Apache Storm [6] was proposed to deal with streaming data.

To better accommodate diverse data processing requirements, the common practice is to co-deploy multiple frameworks in the same cluster and choose the most suitable ones for different applications. Instead of statically partitioning cluster resources for different frameworks, a centralized resource management service is deployed to allocate a certain amount of resources to form a *resource container* at one of the servers

This work was partially supported by National Science Foundation Career Award CNS-1452751, National Science Foundation grant CNS-1552525 and AFOSR grant FA9550-14-1-0160.

to execute a task. Two popular and representative resource management platforms are Hadoop YARN [7] and Apache Mesos [8], which share a similar design with centralized resource allocation and fine-grained resource representation. When the cluster is initially launched, each node declares its resource capacities, e.g., the number of CPU cores and the memory size. Meanwhile, applications from different frameworks send resource requests for their tasks to the centralized resource manager. The resource management tracks the available resources when allocating the containers, and guarantees that the resources occupied by all the containers on a host do not exceed its capacities.

While providing easy management and performance isolation, the existing *exclusive* mode of resource container leads to a potential problem of underutilizing the cluster resources and degrading system performance significantly. For example, a production cluster at Twitter managed by Mesos has reported its aggregated CPU utilization lower than 20% [9] when reservations reach up to 80%. Similarly, Google’s Borg system has reported an aggregated CPU utilization of 25-35% while reserved CPU resources exceed 70% [10]. The major reason of such a low utilization is that the resources (e.g., CPU cores and memory) occupied by a task will not be released until that task is finished. However, tasks from many data processing applications often exhibit fluctuating resource usage patterns. A task may not fully use all the resources throughout its execution. As an example, a reduce task in MapReduce usually has low CPU utilization during its shuffle stage but demands more CPU resources once all intermediate data are received. Another example is an interactive Spark job, where the resource usage of its tasks can be extremely low during a user’s thinking time but significantly increases upon the arrival of a user request.

To solve this problem, we present a new *opportunistic* and *efficient resource allocation* scheme, named **OPERA**, which aims to break the barriers among the encapsulated resource containers by sharing their occupied resources. Our goal is to develop general techniques that can be integrated into a unified resource management framework such that the cluster resources can be shared by multiple data processing paradigms to increase resource utilization and cost efficiency. The main idea of our approach is to leverage the knowledge of actual runtime resource utilizations as well as future resource

availability for task assignments. When a task becomes idle or is not fully utilizing its reserved resources, OPERA re-assigns the idle resources to other pending tasks for execution.

In particular, there are two key problems that we need to consider in the design of this new approach. First, resource usage can dynamically change across time. Second, a server node can be overloaded due to resource over-provisioning, which may incur interference and degrade the performance of all active tasks. In this paper, we present a solution that includes the following features to address the above two issues:

- dynamically monitors the runtime resource utilization;
- classifies the pending tasks to determine if each task is eligible for the new opportunistic resource allocation;
- efficiently assigns the idle resources occupied by the running tasks to other eligible pending tasks, and integrates the new approach with the existing resource allocation; and
- mitigates the severe resource contentions caused by opportunistic resource allocation.

We implement OPERA in Hadoop YARN and evaluate its performance with a set of representative MapReduce and Spark applications. Our experimental results show that our OPERA can significantly reduce the average job execution time and increase the resource (memory and CPU) utilizations.

The organization of this paper is as follows. We present our understandings of task resource usage patterns and the intuition of opportunistic resource allocation in Section II. Section III describes details of our new opportunistic resource allocation approach. The performance of this new approach is evaluated under the workloads mixed with MapReduce and Spark jobs in Section IV. The related work is presented in Section V. We finally give our conclusion in Section VI.

II. MOTIVATION

In the current resource management, a certain amount of resources are *exclusively* allocated to each running task, and will be recycled (i.e., re-assigned) only after the task is completed. This mechanism works well with short-lived, fine-grained tasks that usually process a consistent workload throughout their executions. When the cluster resources are repeatedly assigned to serve this type of tasks, the entire system can consistently keep a high resource utilization. However, if data processing jobs include tasks with long life cycles, the current resource management may not work efficiently. Long tasks usually consist of multiple internal stages and may have different resource usage patterns in each stage. For example, reduce tasks include two stages: data transfer/shuffling and reduce. Network bandwidth is the main resource consumed in the former stage and CPU resources are mainly used in the latter stage. The current framework allocates a fixed amount of resources throughout the life time of the task often leading to low resource utilization. In the above example of reduce tasks, the CPU utilization is low in the first stage of shuffling.

To explore this issue, we conduct experiments in a YARN cluster of 20 slave nodes (8 CPU cores and 16GB memory per node) to better understand the task resource usage patterns. In these experiments, we launch a *TeraSort* job (sorting

a randomly generated 50GB input data) on a MapReduce Hadoop platform and a *pathSim* job (a data mining algorithm to analyze the similarity between authors using the academic paper submission records [11]) on a Spark platform. Fig. 1 shows the measured CPU and memory utilizations during the execution of different tasks. We observe that resource usages of all tasks, especially CPU utilizations, are fluctuating over the time. For example, the CPU utilizations of reduce tasks (see Fig. 1(b)) are extremely low for a long period (e.g., 100s ~ 140s) because reduce tasks are waiting for the output of map tasks. Similarly, the resource usages of a Spark job change across time, see Fig. 1(c) and (f). The CPU utilizations of that Spark job range from 20% to 96%. Meanwhile, as Spark is a memory processing framework, we request 9GB memory for each task of that Spark job. However, we find that the assigned memory (i.e., 9 GB) is not always fully utilized by each task, especially at the beginning of the processing, as shown in Fig. 1(f).

One possible solution to avoid low resource utilizations is to use time series data (such as actual resource demands) for resource allocation. However, this solution is barely practical. It is difficult to accurately predict actual resource usages of tasks running in many frameworks. It can also dramatically increase the complexity of resource allocation even if a precise prior knowledge is available. Therefore, we consider an alternative solution that attempts to improve resource utilizations by opportunistically allocating resources for tasks, i.e., re-assigning the occupied but idle resources to other pending tasks. However, we find that simply assigning the idle resources to a random pending task may raise the risk of severe resource contention. We need to develop a new approach that can gain the benefits of opportunistic scheduling by identifying an appropriate set of idle resources and assigning them to suitable pending tasks.

III. THE DESIGN OF OPERA

In this section, we present a new opportunistic resource allocation approach, named OPERA, which leverages the knowledge of the actual runtime resource utilizations to determine the availability of system resources and dynamically re-assigns idle resources to the waiting tasks. The primary goal of this design is to break the barriers among the encapsulated resource containers and share the reserved resources among different tasks (or containers) such that the overall resource utilization and system throughput can be improved.

A. Sketch of OPERA

As a resource management scheme, OPERA's goal is to assign the available resources in the system to the pending tasks. However, the definition of “available resources” in OPERA is different from that in the traditional systems. They include not only the resources that have not yet been assigned to any tasks, but also the occupied resources that are idle at the runtime. Therefore, OPERA includes two types of resource allocations, *normal resource allocation* and *opportunistic resource allocation*, referring to assigning the

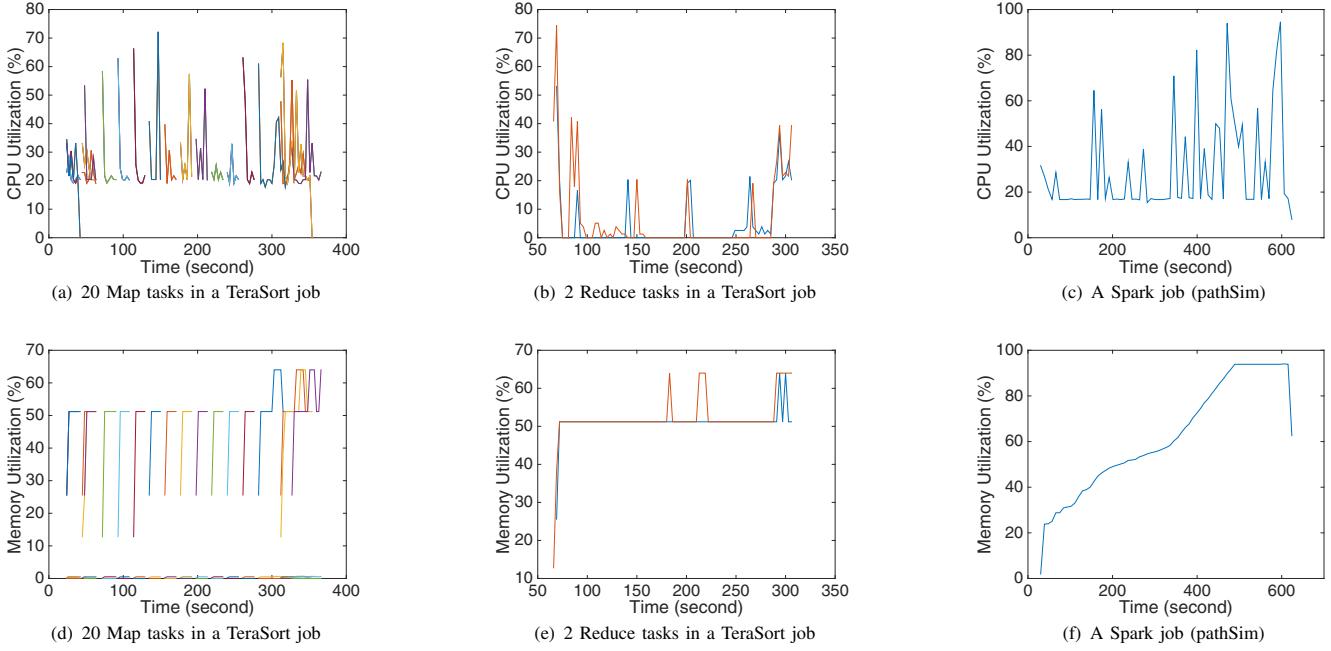


Fig. 1: The CPU and memory utilizations of the tasks from a MapReduce *TeraSort* job and a Spark *pathSim* job.

former and the latter types of available resources, respectively. Then, the basic design of OPERA boils down to two tasks, identifying the available resources under the new definition, and selecting a candidate pending task for execution.

We first define some terms that will be used in our design:

- **Opportunistic/guaranteed available resources:** When assigning *guaranteed available resources* to a task, our OPERA system always guarantees that those resources are available throughout that task's lifetime. On the other hand, if a task is assigned with *opportunistic available resources*, it might lose these resources and get terminated during its execution.
- **Opportunistic/normal tasks:** The tasks that are served with/without opportunistic available resources.

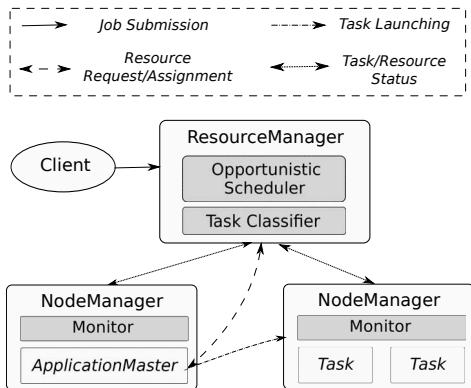


Fig. 2: The architecture of the OPERA-based YARN framework. The modified and newly designed components are marked in grey.

In particular, we develop OPERA on the top of the existing

Hadoop YARN framework. The architecture is illustrated in Fig. 2. We develop the following three major components (the grey parts in Fig. 2).

Task Classifier: The goal of this component is to identify the eligible pending tasks for opportunistic resource allocation. As we discussed in Section II, opportunistic resource allocation is not suitable for all tasks. It is based on resource over-provisioning, and could cause severe resource contentions. In our design, only *short* tasks are eligible for opportunistic resource allocation because longer tasks are more likely to cause resource contentions. However, estimating the execution time of task is challenging in practice. This component is developed to dynamically and accurately classify all the pending tasks into two categories, i.e., short and long.

NodeManager Monitor: This component runs on each cluster node and mainly provides two functions. First, it monitors the dynamic usage of the occupied resources, i.e., the CPU and memory serving the active tasks on the node, and periodically reports the status to the ResourceManager. This function helps the ResourceManager estimate the available resources on each node for opportunistic resource allocation. Accurate reports in a timely fashion are crucial to the system performance. The second function of this component is to mitigate severe resource contentions when the total effective resource utilization is close to or over 100%. Multiple strategies are adopted in this component.

Opportunistic Scheduler: The last component is the core in our system. This component manages both normal and opportunistic resource allocations. When applying opportunistic resource allocation, this component identifies the available resources in the system based on the actual resource usage collected by *NodeManager Monitor*, and allocates them to the

eligible pending tasks determined by *Task Classifier*.

The details of the above components are discussed in the following subsections. Table I lists a summary of notations used in this section.

TABLE I: Notations in this section

Notation	Description
n_i / t	node i / task t
r	a general resource $r \in \{\text{CPU, memory}\}$
$C_i(r)$	resource capacity (r) of n_i
NT_i / OT_i	normal / opportunistic task set on n_i
RT_i	set of all running tasks on n_i , $RT_i = NT_i \cup OT_i$
$D_t(r)$	resource demand (r) of task t
$NU_i(r)$	total resource usage (r) on n_i
$TU_t(r)$	resource usage (r) of task t , $NU_i(r) = \sum_{t \in RT_i} TU_t$
$GA_i(r)$	guaranteed available resource (r) on n_i
$OA_i(r)$	opportunistic available resource (r) on n_i

B. Task Classifier

Task Classifier is expected to classify all the pending tasks into two categories (i.e., short and long), indicating whether they are eligible for the opportunistic resource allocation. Our design basically includes two steps: (1) estimate the execution time of a task; (2) compare it to a threshold to determine the category it belongs to. The classification accuracy is the major concern in our design.

The execution time of a task is hard to be profiled before the execution because it depends on some runtime parameters such as the hardware of the processing cluster. Prior work [12], [13] attempted to use historic statistics of the same type of tasks in the same job to estimate the execution time. We adopt the same approach and extend it to consider the historic information from other applications and even other frameworks. The intuition is that the tasks from the same processing stage of the same type of applications usually have an identical function, process the input data with similar size, and thus have similar execution time. For example, in a MapReduce system, the input data size of each map task is configured by a system parameter. Consider a cluster processing two *wordcount* applications such that one with 10GB input data and the other with 100GB input data. Different input data sizes only yield different numbers of map tasks. All map tasks in these two applications will process similar size of input files and have similar execution time.

Obviously, the historic task execution information is helpful. However, the challenge here is to estimate the execution time when there is no exactly matching historic information. For example, the first batch of map tasks of a new type of MapReduce application and reduce tasks in general MapReduce jobs (usually there are only a few reduce tasks and they may not become effective references to each other). Our solution aims to derive an estimation based on the information from different processing stages, applications, and frameworks.

In OPERA, we adopt the naive Bayes classifier [14] to identify tasks as *short* or *long*. It has been widely used in text classification and spam filtering due to its high accuracy and low overhead on both storage and computation. And we

address the challenges of task classification by presenting a new hierarchy approach that considers the following five properties of each task t .

- F_t : the framework of task t , e.g., MapReduce, Spark, and Storm;
- A_t : the application name of task t , e.g., wordcount and sort;
- S_t : the processing stage that task t belongs to;
- P_t : the progress of the application that task t belongs to;
- D_t : the resource demands of task t .

The first three properties ($\{F_t, A_t, S_t\}$) identify the computation process of each task t , e.g., { MapReduce, WordCount, Map }, and { Spark, PageRank, Stage 3 }. Their values represent the computation details with different granularities. The last two properties are runtime parameters configured by the users. Apparently, with different resource demands (D_t), the execution time of the same task could vary. The progress of the application (P_t) is another implicit factor that may affect the execution time. For example, the user can configure a MapReduce application such that its reduce tasks will start after all its map tasks are finished or the reduce tasks can start when half of the map tasks are completed. In either of these configurations, given different progress values (e.g., all or half of map tasks finish) of the job, the execution time of a reduce task will be different.

We then define the features of each task t as a tuple using Eq. 1.

$$\mathcal{F}_t = \{ \{F_t\}, \{F_t, A_t\}, \{F_t, A_t, S_t\}, \\ \{F_t, A_t, S_t, P_t\}, \{F_t, A_t, S_t, P_t, D_t\} \} \quad (1)$$

We find that combining these task properties together to form such hierarchy features provides more meaningful hints for a better prediction accuracy. In fact, considering each of these properties individually does not provide useful information for classification. For example, the map tasks (same S_t) in different MapReduce applications may yield different execution times; the tasks from the same application “sort” (same A_t) in MapReduce and Spark frameworks may not have the same execution time. However, on the other extreme side, if we classify tasks only based on the historic information from the tasks with the same values of all properties, there will be a lack of information for many tasks and we will miss the correlation between the tasks that share a subset of the properties. Therefore, we decide to combine different task properties in a hierarchical structure in order to explore the correlation between “similar” tasks and confine our estimation in a reasonable scope to accurately classify the tasks.

Once a task t ’s features (\mathcal{F}_t) are determined, we calculate the posterior probability ($P(C|F_t)$) of its category (C_t) using Eq. 2–3 as follows.

$$P(C_t|\mathcal{F}_t) \propto P(C_t) \cdot P(\mathcal{F}_t|C_t), \quad (2)$$

$$P(\mathcal{F}_t|C_t) = \prod_i P(\mathcal{F}_t^i|C_t), \quad (3)$$

where $C_t \in \{\text{short, long}\}$ and \mathcal{F}_t^i represents the i th element of the feature tuple \mathcal{F}_t . Task t is then classified to one of

the two categories which yields a higher posterior probability. Probabilities, e.g., $P(C_t)$, $P(\mathcal{F}_t^i | C_t)$ used in Eq. 2–3, are online learned and updated upon the completion of tasks. We determine the category (short or long) of the finished tasks by checking if their execution times are less than a threshold (e.g., 1 minute) and update all the related probabilities with tasks’ features and category information. There is a special case when an application with all new features is submitted, i.e., no historical information can be referred. In our solution, we opt to conservatively classify the task as a *long* one.

C. NodeManager Monitor

The key idea of our new opportunistic scheduling scheme is to assign idle resources to the pending tasks based on the actual runtime resource usages. Therefore, we develop a monitor module on NodeManagers to (1) keep tracking both CPU and memory usages of the running tasks and sending the collected usage information to ResourceManager through heartbeat messages; and (2) detect and solve performance interferences caused by resource contentions when the resources on a node have been over provisioned and the overall resources occupied by the running tasks exceed the node’s capacity.

Algorithm 1: Node Monitoring

```

Data:  $C_i(r)$ ,  $POLICY$ ,  $BR_i$ ,  $BT_i$ 
Procedure Monitoring()
1   while TRUE do
2      $NU_i(r) \leftarrow 0$ ,  $op \leftarrow$  false,  $c \leftarrow$  “NONE”;
3     foreach  $t$  in  $RT_i$  (the set of running tasks) do
4        $NU_i(r) \leftarrow NU_i(r) + CurrentUsage(r)$ ;
5       if  $t$  is an opportunistic task then
6          $| op \leftarrow$  true;
7       if  $op$  then
8          $| c = CRes(NU_i)$ ;
9       RelieveContention( $c$ ,  $POLICY$ );
10      SLEEP MonitorInterval;
11
12 Procedure CRes( $NU_i$ )
13   if  $NU_i(\text{mem}) > \rho * C_i(\text{mem})$  then
14     return “Memory”;
15   if  $NU_i(\text{CPU}) > \rho * C_i(\text{CPU})$  then
16     return “CPU”;
17   return “NONE”;
18
19 Procedure RelieveContention( $c$ ,  $PO$ )
20   if  $PO = AGGRESSIVE$  and  $c = Memory$  then
21     kill the most recently launched opportunistic task;
22   else if  $PO = NEUTRAL$  and  $c \neq NONE$  then
23     kill the most recently launched opportunistic task;
24   else if  $PO = PRESERVE$  then
25     if  $c \neq NONE$  then
26       kill the most recently launched opportunistic task;
27        $BR_i \leftarrow 2 \cdot BR_i$ ;
28        $BT_i \leftarrow 2 \cdot BT_i$ ;
29        $LastReliefTime = CurrentTime$ ;
30     else
31       if  $CurrentTime - LastReliefTime > BT_i$ 
32         then
33            $| BR_i \leftarrow BR_i/2$ ;
34            $| BT_i \leftarrow BT_i/2$ ;
35            $| LastReliefTime = CurrentTime$ ;

```

Algorithm 1 shows the main process for monitoring resource utilization, detecting and mitigating resource contention on a working node which consists of three modules. In particular, the first module (lines 1–11) periodically collects the CPU and memory usages of the running tasks. $NU_i(\text{CPU})$ and $NU_i(\text{mem})$ represent the CPU and memory usage on node n_i , respectively. We use op to check if there exists any opportunistic task on the node (lines 6–7). If op is false, there will be no resource contention caused by the opportunistic resource allocation. Otherwise, resource contention is possible, and we call the function $CRes$ (line 9) to check the effective resource utilization and return the type of contended resource indicated by variable c . Eventually, the algorithm calls *RelieveContention* function to mitigate the resource contention. The arguments passed to the function are the type of the resource that causes the contention (“CPU”, “Memory”, or “NONE”), and the user-specified policy for handling the contention.

The second module, $CRes$ (lines 12–17), simply compares the resource usage $NU_i(\text{CPU})$ and $NU_i(\text{mem})$ with a pre-defined threshold. If the threshold has been exceeded, the algorithm determines that there exist resource contentions, and reports the type of the contended resource to the main monitoring module. In the algorithm, we set the threshold as $\rho * C_i$, where $C_i(\text{CPU})$ and $C_i(\text{mem})$ are the CPU and memory capacity of node n_i . ρ is an adjusting parameter that can tune the performance of our scheme. By default, we set ρ to 0.95. Note that if both CPU and memory have contentions, this module returns “Memory” as we give memory contention a higher priority to be mitigated.

Finally, in the third module *RelieveContention*, we consider the following three policies to solve the problem of performance interference caused by resource contentions.

- **AGGRESSIVE**: this policy kills the most recently launched opportunistic task only when the monitor detects contention on memory resources;
- **NEUTRAL**: this policy kills the most recently launched opportunistic task under either CPU or memory contention;
- **PRESERVE**: this policy applies the same behaviors as **NEUTRAL**. It further blocks some opportunistic available resources on the node for a period of time;

In all three policies, when resource contentions are detected, the NodeManager attempts to kill the most recently launched opportunistic task to reduce the resource consumption.

AGGRESSIVE policy (lines 19–20) only checks memory contentions. This policy ignores the CPU contention because it is usually less harmful and does not lead to task failures as memory contention does. As opportunistic tasks are relatively short and can release the occupied resources quickly, **AGGRESSIVE** policy tends to aggressively keep these opportunistic tasks running even under CPU contentions for achieving a better overall performance. On the other hand, the drawback of this policy is that the normally reserved resources cannot always be guaranteed especially during the periods of system overloading. In contrast, **NEUTRAL** is a conservative policy (lines 21–22) that kills opportunistic tasks under both

CPU and memory resource contentions. Clearly, this policy can guarantee the reserved resources but might incur frequent task terminations, especially when resource utilizations of the running tasks are oscillating.

To guarantee the reserved resources without killing too many tasks, we further present a PRESERVE policy for contention mitigation, by introducing the concepts of blocked resource (BR_i) and block time (BT_i), see lines 23-33 of Algorithm 1. Besides killing opportunistic tasks, this policy further blocks a certain amount (BR_i) of opportunistic available resources of node n_i for a time window (BT_i). Under the PRESERVE policy, the opportunistic scheduler estimates opportunistic available resources (OA_i) by considering both the actual resource usage of running tasks (TU_t) and the amount of blocked resources (BR_i), as shown in Eq. 4.

$$OA_i(r) = C_i(r) - \sum_{t \in RT_i} TU_t(r) - BR_i(r). \quad (4)$$

The values of BR_i and BT_i are adjusted exponentially in our solution. We double the BR_i and BT_i values to be more conservative if a new resource contention is detected within the current blocking time window. Similarly, the values of BR_i and BT_i are decreased exponentially by a factor of two if no resource contention has been detected in the BT_i window. We also set the minimum/maximum thresholds for both BR_i and BT_i , e.g., in our experiments, we have $0 < BR_i \leq 0.5 \cdot C_i$ and the range of BT_i is between 30 seconds and 90 seconds.

D. Opportunistic Scheduler

Finally, we present the last component in this subsection. As discussed in Section II, the ResourceManager under the current YARN framework considers each resource container exclusively allocated for a single task. When assigning resources to a pending task, the ResourceManager checks the *available resources* on each node as follows,

$$C_i(r) - \sum_{t \in RT_i} D_t(r). \quad (5)$$

However, in practice, the tasks do not always fully utilize their assigned resources during their executions. The traditional resource allocation often leads to a low resource utilization.

To address this issue, we develop the opportunistic scheduler, which considers both *guaranteed available resources* and *opportunistic available resources*. The key difference between these two types of resource allocations is the calculation of the available resources. The guaranteed available resources ($GA_i(r)$) are defined in Eq. 6, which equal to the differences between the resource capacities and the total resource demands of the normal tasks on node n_i . When calculating the opportunistic available resources ($OA_i(r)$), we consider the runtime resource usages of the running tasks rather than their resource demands, see Eq. 7.

$$GA_i(r) = C_i(r) - \sum_{t \in NT_i} D_t(r), \quad (6)$$

$$OA_i(r) = C_i(r) - \sum_{t \in RT_i} TU_t(r), \quad (7)$$

where $C_i(r)$ represents the capacity of resource r of node n_i and $D_t(r)$ and $TU_t(r)$ represent task t 's resource demand and resource usage, respectively.

As discussed in Section III-A, the NodeManager of each working node periodically sends heartbeat messages to the ResourceManager, which includes the node's health status and runtime resource utilizations of each running task. Once receiving a heartbeat message from one working node, our system updates the guaranteed/opportunistic available resources of that node using Eq. 6-7. Then, a pending task is chosen for assignment according to the user defined scheduling policy, e.g., Fair. When allocating available resources, the opportunistic scheduler always first tries to assign guaranteed available resources, i.e., normal resource assignment. If the resource demand of the task cannot be fulfilled, the scheduler then attempts to allocate opportunistic available resources. If the task is not eligible for opportunistic scheduling or still cannot fit into the opportunistic available resources on the node, the scheduler reserves this task on the node and stops the assignment process until receiving the next heartbeat message. Only one task can be reserved on each node, and the reserved task has a higher priority in the resource allocation.

IV. EVALUATION

We implement the proposed opportunistic resource allocation scheme, OPERA, on Hadoop YARN v2.5. Specifically, we modify the scheduler component in the ResourceManager of YARN (on top of the Fair scheduler) to include a task classifier for separating the opportunistic task assignment from the normal task assignment. We note that OPERA can also be integrated with any other scheduling algorithms. In the NodeManager, we enable the mechanisms of runtime resource monitoring/reporting as well as contention detection and integrate these mechanisms in the ContainerMonitor component. The communication protocols and messages among the ResourceManager, NodeManagers, and ApplicationMasters are also modified to convey the actual resource usage and the assignment type (i.e., normal or opportunistic) information of tasks. We evaluate OPERA in a real YARN cluster with different data processing workloads which include mixed sets of representative MapReduce and Spark jobs.

A. Experiment Settings

We conduct our experiments in a YARN cluster which is deployed in a cloud environment provided by CloudLab [15]. This YARN cluster is configured with one master node and 20 working nodes, each of which has 8 physical cores. We configure 2 virtual cores for each physical core such that there are 16 vcores in total on each working node. Among those 16 vcores, we use one vcore for NodeManager and the HDFS usage, and the remaining 15 vcores for running cluster computing applications. Each node is configured with memory capacity of 12 GB. Thus, the total resource capacity of this cluster is $<300\text{vcores}, 240\text{GB}>$ for CPU vcores and memory.

Four benchmarks are considered in our experiments. (1) *PathSim*, a Spark application [11], computes the meta path based on the similarity among academic paper authors. The input data contains 1.2 million paper submission records. (2) *Terasort*, a MapReduce application, sorts the input records. We use 50GB input data generated by *teragen*. (3) *Wordcount*, a MapReduce application, counts the occurrences of each word in the input files. Each input file with 50GB data is generated through *randomTextWriter*. (4) *PiEstimate*, a MapReduce application, estimates the value of π using the quasi-Monte Carlo method. Each task processes 300 million data points.

Table II shows the configurations of each application, including task numbers and task resource requests. By default, we configure each task's resource request according to their actual resource usage. The CPU demand of a task is equal to 75% of its peak CPU usage and the memory requirement is set to that task's maximum memory usage. These applications thus have various resource requirements. For example, the tasks from Spark applications are memory intensive while MapReduce tasks are mainly CPU intensive.

TABLE II: Task Configurations of Applications

Framework	Application	Task Num.	Task Resource Req.
Spark	<i>pathSim</i>	10 executors	$< 4\text{vcores}, 9\text{GB} >$
MapReduce	<i>terasort</i>	374 mappers 100 reducers	$< 4\text{vcores}, 1\text{GB} >$ $< 2\text{vcores}, 1\text{GB} >$
	<i>wordcount</i>	414 mappers 50 reducers	$< 4\text{vcores}, 1\text{GB} >$ $< 2\text{vcores}, 1\text{GB} >$
	<i>piEstimate</i>	500 mappers 1 reducers	$< 3\text{vcores}, 1\text{GB} >$ $< 2\text{vcores}, 1\text{GB} >$

B. Experiment Results

In our experiments, we set the results under the original YARN framework with the Fair scheduler as a baseline for comparison. The major performance metrics we consider for evaluating OPERA include resource utilizations and job execution times.

1) *Workloads with MapReduce jobs only*: In the first set of experiments, we generate a workload of 6 MapReduce jobs, 2 jobs from each MapReduce application listed in Table II. We execute the same set of 6 jobs with the traditional resource allocation under the Fair scheduler, and with our OPERA system with three mitigation policies, i.e., AGGRESSIVE, NEUTRAL, and PRESERVE.

Table III first shows the prediction accuracy of our task classifier which adopts the hierarchy feature-based naive Bayes classification algorithm as described in Section III-B. In this table, the “Fact” column shows the total number of actual short (resp. long) tasks, while the “Classification” column presents the predicted results, i.e., the number of tasks that are predicted to be short and long, and the predication accuracy.

TABLE III: Task classification.

	Fact	Classification		
		Short	Long	Pred. Accuracy
<i>Short</i>	2777	2417	360	87.0%
<i>Long</i>	107	2	105	98.1%

We observe that our task classifier is able to accurately categorize most of short and long tasks with high prediction accuracy ratios, i.e., 87% and 98%, respectively. More importantly, our classifier successfully avoids the false positives (i.e., predicting a long task as short) that can incur severe resource contention and other harmful consequences. As shown in Table III, only 1.9% (i.e., 2 out of 107) of the long tasks are classified as short ones. On the other hand, we notice that the false negative (i.e., predicting a short task as long) ratio is slightly high, i.e., 13%. However, it is still in a low range and only prevents us from assigning opportunistic available resources to those tasks, which in general does not degrade the overall performance.

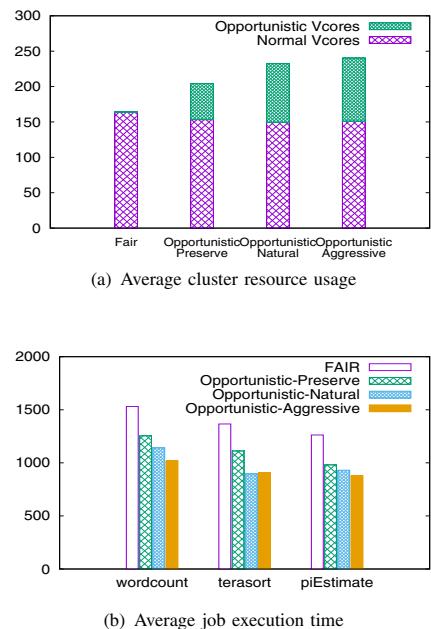


Fig. 3: The performance in the experiments with 6 MapReduce jobs

Fig. 3 presents average cluster resource usages (i.e., the number of used vcores) and average job execution times (in seconds). We observe that our OPERA scheduler is able to more efficiently use cluster resources as shown in Fig. 3(a). As mentioned in Section IV-A, the total CPU capacity in our YARN cluster is 300 vcores, i.e., $15 \text{ vcores} \cdot 20 \text{ nodes}$. We can see that the original YARN system with Fair only uses about 50% (164) of vcores in average. OPERA increases the number of actually used vcores up to 240 (e.g., with the AGGRESSIVE scheme), which accelerates the overall processing by using more CPU resources. As a result, compared to Fair, our OPERA significantly improves the average job execution times for each MapReduce application, see Fig. 3(b).

To better understand how OPERA and Fair work, we present the CPU vcore usages, the number of normal tasks, and the number of opportunistic tasks across time on a single cluster node in Fig. 4 and Fig. 5. Obviously, under the Fair policy, the number of vcores that are actually used is low and fluctuating across time, see Fig. 4(a). Moreover, the number of normal tasks does not change much under Fair (see Fig. 5(a)),

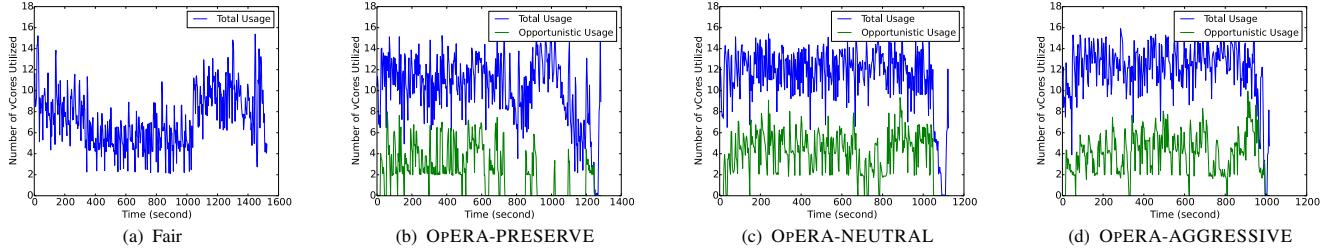


Fig. 4: Runtime CPU usages on a single cluster node under the workload with 6 MapReduce jobs

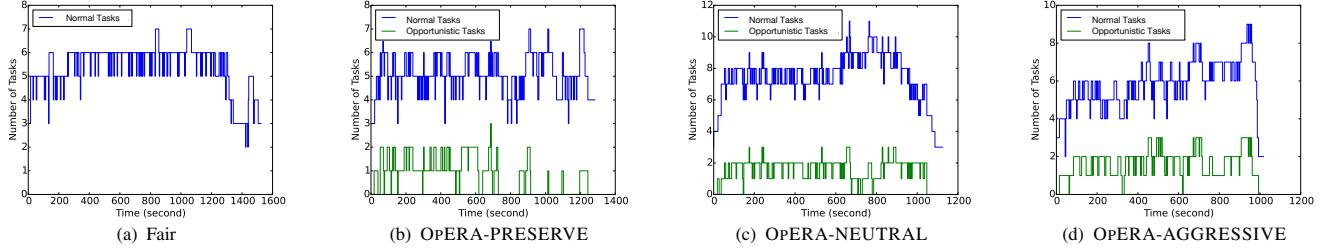


Fig. 5: The numbers of normal tasks and opportunistic tasks on a single cluster node under the workload with 6 MapReduce jobs

which further indicates that these tasks running under the traditional resource allocation with the Fair scheduler yield varying CPU usage patterns. On the other hand, through the opportunistic resource allocation, the system resources (e.g., CPU vcores) are better utilized because more tasks are scheduled to run in a node when we detect underutilized resources on that node, as shown in Fig. 4(b)-(d). Particularly, OPERA with AGGRESSIVE or NEUTRAL always keeps CPU resources fully utilized (i.e., around 15 vcores in use per node), see plots (c) and (d) in Fig. 4.

2) *Workloads with MapReduce and Spark jobs:* In the second set of experiments, we launch two Spark jobs, i.e., *pathSim* [11], together with 6 MapReduce jobs that are the same as we have in the first set of experiments. Here, each Spark executor occupies 10GB (i.e., 9GB executor memory request and 1GB overhead) memory on a cluster node.

Fig. 6 shows the experimental results, including the average job execution time and the average vcore usage under different scheduling policies. Obviously, all MapReduce jobs receive a significant performance improvement under OPERA. As shown in Fig. 6(a), the average job execution time of all MapReduce jobs is reduced by 25.5%, 29.3%, and 36.1% under PRESERVE, NEUTRAL, and AGGRESSIVE, respectively. On the other hand, the two Spark jobs (i.e., *pathSim*) do not benefit from our opportunistic resource allocation. The reason is that all tasks in *pathSim* are launched together in a single wave. The parallelism of these Spark jobs thus cannot be further improved through opportunistic resource allocation. Moreover, the performance of two Spark jobs becomes slightly worse under our opportunistic resource allocation due to the resource contention caused by other opportunistically scheduled MapReduce tasks.

V. RELATED WORK

Improving resource efficiency and throughput of cluster computing platforms was extensively studied in recent years.

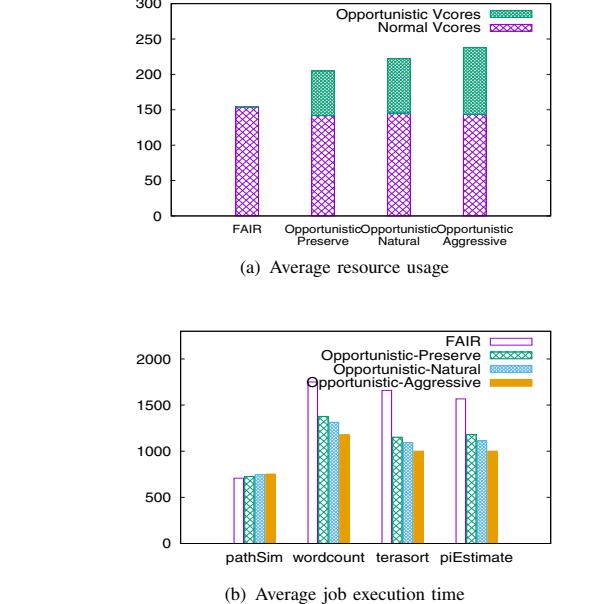


Fig. 6: The experiment with MapReduce and Spark jobs.

Our previous works [16], [17] focus on the first generation Hadoop system which adopts coarse-grained resource management. For fine-grained resource management, we proposed a scheduler HaSTE in Hadoop YARN [18] that improves resource utilization using more efficient task packing according to diverse resource requirements of tasks on different resource types and dependencies between tasks. However, HaSTE only considers the task dependencies in the MapReduce framework and assigns resources according to task requests without considering real time resource usages. DynMR [19] presents that reduce tasks in the MapReduce framework bundle multiple phases and have changing resource utilization, and proposes

to assemble multiple reduce tasks into a progressive queue for management and backfill map tasks to fully utilize system resources. Their work is closely bounded with the MapReduce framework, and involves complex task management that cannot be easily extended to other frameworks. Quasar [9] designs resource efficient and QoS-aware cluster management. Classification techniques are used to find appropriate resource allocations to applications in order to fulfill their QoS requirements and maximize system resource utilization. Resource assignment in their work is to assign one or multiple nodes to the application, which is different from task assignment in cluster computing. The previous studies mainly address the inefficient resource utilization caused by the gap between user specified application resource requirements and actual resource usages of applications. While, our work mainly addresses the issue of resource underutilization that is caused by the fluctuating resource usage patterns of tasks.

Resource utilization is of greater importance in large enterprise data centers since increasing utilization by few percentages can save a lot in a large-scale cluster. Recent published works reveal some technique details of Google's Borg [20] and Microsoft's Apollo [21] systems. Similar to our solution, they both consider the runtime resource utilization. Borg classifies jobs into the categories of high priority and low priority. If high priority tasks are not using all their reserved resources, resource manager can reclaim these resources and assign them to low priority tasks. Apollo starts opportunistic scheduling after all available resource tokens have been assigned to regular tasks. Our OPERA system, however, includes enhanced schemes in the following aspects. First, instead of using user-defined task priorities or scheduling the fixed amount of opportunistic tasks, our solution dynamically recognizes the appropriate pending tasks for the opportunistic resource allocation. As a result, the interference introduced by opportunistic scheduling and the penalty of killing unfinished opportunistic tasks can be minimized under our proposed approach. In addition, our approach uses three resource release schemes that consider different degrees of aggressiveness of opportunistic scheduling.

VI. CONCLUSION

In this paper, we developed a novel resource management scheme, named OPERA, to enable the sharing of occupied resources among different tasks (or resource containers). To meet this goal, OPERA leverages the knowledge of actual runtime resource utilizations to detect underutilized resources and opportunistically re-assigns these resources to other eligible pending tasks (i.e., the expected short ones). By this way, we guarantee that performance interference can be minimized and killing opportunistically launched tasks does not lead to a significant waste of work. We implemented OPERA on the top of Hadoop YARN v2.5 and evaluated our proposed scheduler in a cloud environment provided by CloudLab. Diverse workloads mixed with MapReduce and Spark jobs have been produced to evaluate OPERA under different scenarios. The experimental results show that our OPERA is able to achieve

up to 39.8% reduction in average job execution time and 30% increase in resource utilizations.

REFERENCES

- [1] J. Dean and S. Ghemawat, "Mapreduce: simplified data processing on large clusters," *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- [2] T. White, *Hadoop: The definitive guide*. O'Reilly Media, Inc., 2012.
- [3] (2014) Hadoop Users. [Online]. Available: <https://wiki.apache.org/hadoop/PoweredBy>
- [4] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica, "Spark: cluster computing with working sets," in *Proceedings of the 2nd USENIX conference on Hot topics in cloud computing*, 2010, pp. 10–10.
- [5] M. Zaharia, M. Chowdhury, T. Das, A. Dave, J. Ma, M. McCauley, M. J. Franklin, S. Shenker, and I. Stoica, "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*. USENIX Association, 2012, pp. 2–2.
- [6] N. Marz, "A storm is coming: more details and plans for release," *Twitter Engineering*, vol. 42, 2011.
- [7] V. K. Vaipallapalli, A. C. Murthy, C. Douglas, S. Agarwal, M. Konar, R. Evans, T. Graves, J. Lowe, H. Shah, S. Seth *et al.*, "Apache hadoop yarn: Yet another resource negotiator," in *Proceedings of the 4th annual Symposium on Cloud Computing*. ACM, 2013, p. 5.
- [8] B. Hindman, A. Konwinski, M. Zaharia, A. Ghodsi, A. D. Joseph, R. H. Katz, S. Shenker, and I. Stoica, "Mesos: A platform for fine-grained resource sharing in the data center," in *NSDI*, vol. 11, 2011, pp. 22–22.
- [9] C. Delimitrou and C. Kozyrakis, "Quasar: Resource-efficient and qos-aware cluster management," *ACM SIGPLAN Notices*, vol. 49, no. 4, pp. 127–144, 2014.
- [10] C. Reiss, A. Tumanov, G. R. Ganger, R. H. Katz, and M. A. Kozuch, "Heterogeneity and dynamicity of clouds at scale: Google trace analysis," in *Proceedings of the Third ACM Symposium on Cloud Computing*. ACM, 2012, p. 7.
- [11] Y. Sun, J. Han, X. Yan, P. S. Yu, and T. Wu, "Pathsim: Meta path-based top-k similarity search in heterogeneous information networks," *VLDL11*, 2011.
- [12] A. Verma, Ludmila Cherkasova, and R. H. Campbell, "Aria: Automatic resource inference and allocation for mapreduce environments," in *ICAC'11*, 2011, pp. 235–244.
- [13] J. Polo, D. Carrera, Y. Becerra, J. Torres, E. Ayguadé, M. Steinder, and I. Whalley, "Performance-driven task co-scheduling for mapreduce environments," in *Network Operations and Management Symposium (NOMS), 2010 IEEE*. IEEE, 2010, pp. 373–380.
- [14] S. J. Russell, P. Norvig, and S. Chakrabarti, "Artificial intelligence: a modern approach."
- [15] Cloudlab. [Online]. Available: <http://cloudlab.us/>
- [16] Y. Yao, J. Wang, B. Sheng, C. Tan, and N. Mi, "Self-adjusting slot configurations for homogeneous and heterogeneous hadoop clusters," *Cloud Computing, IEEE Transactions on*, no. 99, p. 1, 2015.
- [17] J. Wang, Y. Yao, Y. Mao, B. Sheng, and N. Mi, "Fresh: Fair and efficient slot configuration and scheduling for hadoop clusters," in *Cloud Computing (CLOUD)*, 2014.
- [18] Y. Yao, J. Wang, B. Sheng, J. Lin, and N. Mi, "Haste: Hadoop yarn scheduling based on task-dependency and resource-demand," in *Cloud Computing (CLOUD), 2014 IEEE 7th International Conference on*. IEEE, 2014, pp. 184–191.
- [19] J. Tan, A. Chin, Z. Z. Hu, Y. Hu, S. Meng, X. Meng, and L. Zhang, "Dynmr: Dynamic mapreduce with reducetask interleaving and maptask backfilling," in *Proceedings of the Ninth European Conference on Computer Systems*. ACM, 2014, p. 2.
- [20] A. Verma, L. Pedrosa, M. Korupolu, D. Oppenheimer, E. Tune, and J. Wilkes, "Large-scale cluster management at google with borg," in *Proceedings of the Tenth European Conference on Computer Systems*. ACM, 2015, p. 18.
- [21] E. Boutin, J. Ekanayake, W. Lin, B. Shi, J. Zhou, Z. Qian, M. Wu, and L. Zhou, "Apollo: scalable and coordinated scheduling for cloud-scale computing," in *Proceedings of the 11th USENIX conference on Operating Systems Design and Implementation*. USENIX Association, 2014, pp. 285–300.