

Using a Tunable Knob for Reducing Makespan of MapReduce Jobs in a Hadoop Cluster

Yi Yao
Northeastern University
Boston, MA, USA
yyao@ece.neu.edu

Jiayin Wang Bo Sheng
University of Massachusetts Boston
Boston, MA, USA
{jane,shengbo}@cs.umb.edu

Ningfang Mi
Northeastern University
Boston, MA, USA
ningfang@ece.neu.edu

Abstract—The MapReduce framework and its open source implementation Hadoop have become the defacto platform for scalable analysis on large data sets in recent years. One of the primary concerns in Hadoop is how to minimize the completion length (i.e., makespan) of a set of MapReduce jobs. The current Hadoop only allows static slot configuration, i.e., fixed numbers of map slots and reduce slots throughout the lifetime of a cluster. However, we found that such a static configuration may lead to low system resource utilizations as well as long completion length. Motivated by this, we propose a simple yet effective scheme which uses *slot ratio* between map and reduce tasks as a tunable knob for reducing the makespan of a given set. By leveraging the workload information of recently completed jobs, our scheme dynamically allocates resources (or slots) to map and reduce tasks. We implemented the presented scheme in Hadoop V0.20.2 and evaluated it with representative MapReduce benchmarks at Amazon EC2. The experimental results demonstrate the effectiveness and robustness of our scheme under both simple workloads and more complex mixed workloads.

I. INTRODUCTION

In recent years, MapReduce [1] has become the leading paradigm for parallel big data processing. Its open source implementation Apache Hadoop [2] has also emerged as a popular platform for daily data processing and information analysis. With the rise of cloud computing, MapReduce is no longer just for internal data process in big companies. It is now convenient for a regular user to launch a MapReduce cluster on the cloud, e.g., AWS MapReduce, for data-intensive applications. When more and more applications are adopting the MapReduce framework, how to improve the performance of a MapReduce cluster becomes a focus of research and development. Both academia and industry have put tremendous efforts on job scheduling, resource management, and Hadoop applications [3]–[11]. As a complex system, Hadoop is configured with a large set of system parameters. While it provides the flexibility to customize the cluster for different applications, it is challenging for users to understand and set the optimal values for those parameters. In this paper, we aim to develop algorithms for adjusting a basic system parameter to improve the performance of makespan of a batch of jobs.

A classic Hadoop cluster includes a single master node and multiple slave nodes. The master node runs the *JobTracker*

routine which is responsible for scheduling jobs and coordinating the execution of tasks of each job. Each slave node runs the *TaskTracker* daemon for hosting the execution of MapReduce jobs. The concept of “slot” is used to indicate the capacity of accommodating tasks on each node. In a Hadoop system, a slot is assigned as a map slot or a reduce slot serving map tasks or reduce tasks, respectively. At any given time, only one task can be running per slot. The number of available slots per node indeed provides the maximum degree of parallelization in Hadoop. Our experiments have shown that the slot configuration has a significant impact on system performance. The Hadoop framework, however, uses fixed numbers of map slots and reduce slots at each node as the default setting throughout the lifetime of a cluster. The values in this fixed configuration are usually heuristic numbers without considering job characteristics. Therefore, this static setting is not well optimized and may hinder the performance improvement of the entire cluster.

In this work, we propose and implement a new mechanism TuMM to dynamically allocate slots for map and reduce tasks. The primary goal of the new mechanism is to improve the completion time (i.e., makespan) of a batch of MapReduce jobs while retain the simplicity in implementation and management of the slot-based Hadoop design. The key idea of TuMM is to automate the slot assignment ratio between map and reduce tasks in a cluster as a tunable knob for reducing the makespan of MapReduce jobs. The Workload Monitor (WM) and the Slot Assigner (SA) are the two major components introduced by TuMM. The WM that resides in the *JobTracker* periodically collects the execution time information of recently finished tasks and estimates the present map and reduce workloads in the cluster. The SA module takes the estimation to decide and adjust the slot ratio between map and reduce tasks for each slave node. With TuMM, the map and reduce phases of jobs could be better pipelined under priority based schedulers, and thus the makespan is reduced.

The rest of the paper is organized as follows. We explain the motivation of our work through some experimental examples in Section II. We formulate the problem and derive the optimal setting for static slot configuration in Section III. The design details of the dynamic mechanism are presented in Section IV. Section V provides the experimental evaluation of the proposed scheme. Section VI describes the related work

of this article. We conclude in Section VII.

II. MOTIVATION

Currently, the Hadoop framework uses fixed numbers of map slots and reduce slots on each node throughout the lifetime of a cluster. However, such a fixed slot configuration may lead to low resource utilizations and poor performance especially when the system is processing varying workloads. We here use two simple cases to exemplify this deficiency. In each case, three jobs are submitted to a Hadoop cluster with 4 slave nodes and each slave node has 4 available slots. Details of the experimental setup are introduced in Section V. To illustrate the impact of resource assignments, we also consider different static settings for map and reduce slots on a slave node. For example, when the slot ratio is equal to 1:3, we have 1 map slot and 3 reduce slots available per node. We then measure the overall lengths (i.e., makespans) for processing a batch of jobs, which are shown in Fig. 1.

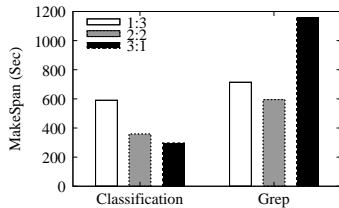


Fig. 1. The makespans of jobs under case 1 (i.e., *Classification*) and case 2 (i.e., *Grep*). The map and reduce slot ratios on each slave node are set to 1:3, 2:2, and 3:1.

Case 1: We first submit three *Classification* jobs to process a 10 GB movie rating data set. We observe that the makespan of these jobs is varying under different slot ratio settings and the best performance (i.e., shortest makespan) is achieved when each slave node has three map slots and one reduce slot, see the left column of Fig. 1.

To interpret this effect, we further plot the execution times of each task in Fig. 2. Clearly, *Classification* is a map-intensive application; for example, when we equally distribute resources (or slots) between map and reduce tasks, i.e., with the slot ratio of 2:2, the length of a map phase is longer than that of a reduce phase, see Fig. 2(a). It follows that each job’s reduce phase (including shuffle operations and reduce operations) overlaps with its map phase for a long period. However, as the reduce operations can only start after the end of the map phase, the occupied reduce slots stay in shuffle for a long period, mainly waiting for the outputs from the map tasks. Consequently, system resources are underutilized.

For example, we tracked the CPU utilizations of each task in a slave node every 5 seconds and Table I shows part of the records in one of such overlapping periods. At each moment, the overall CPU utilization (i.e., the summation of CPU utilizations of the four tasks) is much less than 400%, for a node with 4 cores. We then notice that when we assign more slots to map tasks, e.g., with the slot ratio of 3:1, each job experiences a shorter map phase and most of its reduce phase overlaps with the following job’s map phase, see Fig. 2(b). The average CPU utilization is also increased

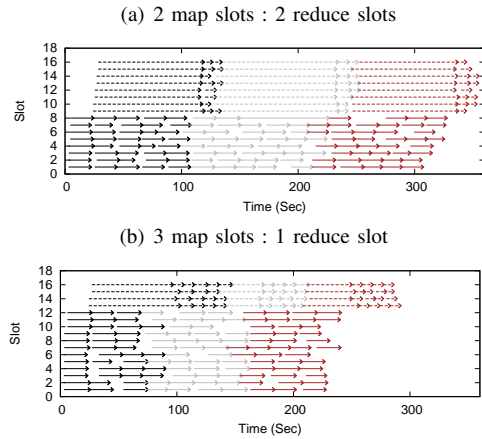


Fig. 2. Task execution times of three *Classification* jobs under different static slot configurations, where each node has (a) 2 map slots and 2 reduce slots, and (b) 3 map slots and 1 reduce slot. Each arrowed line represents the execution of one task, and the solid (resp. dashed) ones represent map (resp. reduce) tasks. The first wave in each job’s reduce phase represents the shuffle operations. In addition, we use three different colors to discriminate the three jobs.

by 20% compare to those under the the slot ratio of 2:2. It implies that for map-intensive jobs like *Classification*, one should assign more resources (slots) to map tasks in order to improve the performance in terms of makespan.

TABLE I
REAL TIME CPU UTILIZATIONS OF EACH TASK ON A SLAVE NODE IN THE OVERLAPPING TIME PERIOD OF A JOB’S MAP AND REDUCE PHASES. THE SLOT RATIO PER NODE IS 2:2.

Time(sec)	ProcessId / TaskType			
	3522 / map	3564 / map	3438 / reduce	3397 / reduce
1	147%	109%	26%	0%
6	103%	93%	0%	4%
11	93%	99%	8%	0%
16	100%	100%	0%	0%
21	97%	103%	0%	0%

Case 2: In this case, we turn to consider reduce-intensive applications by submitting three *Grep* jobs to scan the 10 GB movie rating data. Similar to case 1, we also investigate three static slot configurations.

First, we observe that each job takes longer time to process its reduce phase than its map phase when we have 2 map and 2 reduce slots per node, see Fig. 3(a). Based on the observation in case 1, we expect a reduced makespan when assigning more slots to reduce tasks, e.g., with the slot ratio of 1:3. However, the experimental results show that the makespan under this slot ratio setting (1:3) becomes even longer than that under the setting of 2:2, see the right column of Fig. 1. We then look closely at the corresponding task execution times, see Fig. 3(b). We find that the reduce tasks indeed have excess slots such that the reduce phase of each job starts too early and wastes time waiting for the output from its map phase. In fact, a good slot ratio should be set between 2:2 and 1:3 to enable each job’s reduce phase to fully overlap with the following job’s map phase rather than its own map phase.

In summary, in order to reduce the makespan of a batch of jobs, more resources (or slots) should be assigned to map (resp. reduce) tasks if we have map (resp. reduce) intensive

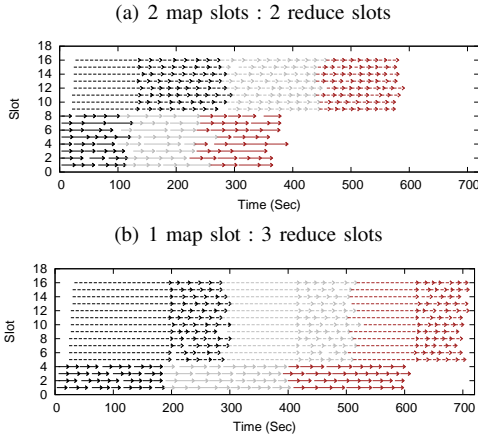


Fig. 3. Task execution times of a batch of *Grep* jobs under different static slot configurations, where each node has (a) 2 map slots and 2 reduce slots, and (b) 1 map slot and 3 reduce slots.

jobs. On the other hand, a simple adjustment in such slot configurations is not enough. An effective approach should tune the slot assignments such that the execution times of map and reduce phases can be well balanced and the makespan of a given set can be reduced to the end.

III. SYSTEM MODEL AND STATIC SLOT CONFIGURATION

In this section, we present the system model we considered and formulate the problem. In addition, we analyze the default static slot configuration in Hadoop and present an algorithm to derive the best configuration.

A. Problem Formulation

In our problem setting, we consider that a Hadoop cluster consisting of k nodes has received a batch of n jobs for processing. We use J to represent the set of jobs, $J = \{j_1, j_2, \dots, j_n\}$. Each job j_i is configured with $n_m(i)$ map tasks and $n_r(i)$ reduce tasks. Let $st(i)$ and $ft(i)$ indicate the start time and the finish time of job j_i , respectively. In addition, we assume the Hadoop system sets totally S slots on all the nodes in the cluster. Let s_m and s_r be the number of map slots and reduce slots, respectively. We then have $S = s_m + s_r$. In this paper, our objective is to develop an algorithm to dynamically tune the parameters of s_m and s_r , given a fixed value of S , in order to minimize the makespan of the given batch of jobs, i.e., $\text{minimize}\{\max\{ft(i), \forall i \in [1, n]\}\}$.

In a Hadoop system, the makespan of multiple jobs also depends on the job scheduling algorithm which is coupled with our solution of allocating the map and reduce slots on each node. In this paper, we assume that a Hadoop cluster uses the default FIFO (First-In-First-Out) job scheduler because of the following two reasons. First, given n jobs waiting for service, the performance of FIFO is no worse than other schedulers in terms of makespan. In the example of ‘‘Case 2’’ mentioned in Section II, the makespan under FIFO is 594 sec while Fair, another alternative scheduler in Hadoop, consumes 772 sec to finish jobs. Second, using FIFO simplifies the performance analysis because generally speaking, there are fewer concurrently running jobs at any time. Usually two jobs, with one in map phase and the other in reduce phase.

Furthermore, we use execution time to represent the workload of each job. As a MapReduce job is composed of two phases, we define $w_m(i)$ and $w_r(i)$ as the workload of map phase and reduce phase in job j_i , respectively. We have developed solutions with and without the prior knowledge of the workload and we will discuss how to obtain this information later.

B. Static Slot Configuration with Workload Information

First, we consider the scenario that the workload of a job is available and present the algorithm for static slot configuration which is default in a Hadoop system. Basically, the Hadoop cluster preset the values of s_m and s_r under the constraint of $S = s_m + s_r$ before executing the batch of jobs, and the slot assignment will not be changed during the entire process. We have developed the following Algorithm 1 to derive the optimal values of s_m and s_r .

Our algorithm and analysis are based on an assumption that the workload of map or reduce phase is inversely proportional to the number of slots assigned to the phase. Given s_m and s_r , the map (resp. reduce) phase of j_i needs $\frac{n_m(i)}{s_m}$ (resp. $\frac{n_r(i)}{s_r}$) rounds to finish. In each round, s_m map tasks or s_r reduce tasks are processed in parallel and the time consumed is equal to the execution time of one map or one reduce task. Let $\bar{t}_m(i)$ and $\bar{t}_r(i)$ be the average execution time for a map task and a reduce task, respectively. The workloads of map and reduce phases are defined as

$$w_m(i) = n_m(i) \cdot \bar{t}_m(i), w_r(i) = n_r(i) \cdot \bar{t}_r(i). \quad (1)$$

Algorithm 1 can derive the best static setting of s_m and s_r given the workload information. The outer loop (lines 1–10) in the algorithm enumerates the value of s_m and s_r (i.e., $S - s_m$). For each setting of s_m and s_r , the algorithm first calculates the workload ($w_m(i)$ and $w_r(i)$) for each job j_i in lines 3–5. The second inner loop (lines 6–8) is to calculate the finish time of each job. Under the FIFO policy, there are at most two concurrently running jobs in the Hadoop cluster. Each job’s map or reduce phase cannot start before the precedent job’s map or reduce phase is finished (we assume here that all jobs have more tasks than the slots number in system for the simplicity of discussion). More specifically, the start time of map tasks of job j_i , i.e., $st(i)$, is the finish time of j_{i-1} ’s map phase, i.e., $st(i) = st(i-1) + \frac{w_m(i-1)}{s_m}$. Additionally, the start time of j_i ’s reduce phase should be no earlier than both the finish time of j_i ’s map phase and the finish time of j_{i-1} ’s reduce phase. Therefore, the finish time of j_i is $ft(i) = \max(st(i) + \frac{w_m(i)}{s_m}, ft(i-1)) + \frac{w_r(i)}{s_r}$. Finally, the variables Opt_SM and Opt_MS keep track of the optimal value of s_m and the corresponding makespan (lines 9–10), and the algorithm returns Opt_SM and $S - Opt_SM$ as the values for s_m and s_r at the end. The time complexity of the algorithm is $O(S \cdot n)$.

IV. DYNAMIC SLOT CONFIGURATION

As discussed in Section II, the default Hadoop cluster uses static slot configuration and does not perform well for varying

Algorithm 1 Static Slot Configuration

```

1: for  $s_m = 1$  to  $S$  do
2:    $s_r = S - s_m$ 
3:   for  $i = 1$  to  $n$  do
4:      $w_m(i) = n_m(i) \cdot \bar{t}_m(i)$ 
5:      $w_r(i) = n_r(i) \cdot \bar{t}_r(i)$ 
6:   for  $i = 1$  to  $n$  do
7:      $st(i) = st(i-1) + \frac{w_m(i-1)}{s_m}$ 
8:      $ft(i) = \max(st(i) + \frac{w_m(i)}{s_m}, ft(i-1)) + \frac{w_r(i)}{s_r}$ 
9:   if  $ft(n) < Opt\_MS$  then
10:     $Opt\_MS = ft(n)$ ;  $Opt\_SM = s_m$ 
11: return  $Opt\_SM$  and  $S - Opt\_SM$ 

```

workloads. The inappropriate setting of s_m and s_r may lead to extra overhead because of the following two cases:

- (1) if job j_i 's map phase is completed later than job j_{i-1} 's reduce phase, then the reduce slots will be idle for the interval period of $(st(i) + w_m(i)) - ft(i-1)$, see Fig. 4(a);
- (2) if job j_i 's map phase is completed earlier than the job j_{i-1} 's reduce phase, then j_i 's reduce tasks have to wait for a period of $ft(i-1) - (st(i) + w_m(i))$ until reduce slots are released by j_{i-1} , see Fig. 4(b).

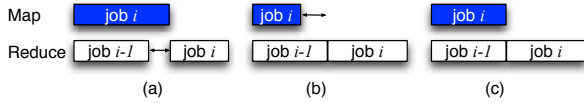


Fig. 4. Illustration of aligning the map and reduce phases. (a) and (b) are the two undesired cases mentioned above, and our goal is to achieve (c).

In this section, we present our solutions that dynamically allocate the slots to map and reduce tasks during the execution of jobs. The architecture of our design is shown in Fig. 5. In dynamic slot configuration, when one slot becomes available upon the completion of a map or reduce task, the Hadoop system will re-assign a map or reduce task to the slot based on the current optimal values of s_m and s_r . There are totally $\sum_{i \in [1, n]} (n_m(i) + n_r(i))$ tasks and at the end of each task, Hadoop needs to decide the role of the available slot (either a map slot or a reduce slot). In this setting, therefore, we cannot enumerate all the possible values of s_m and s_r (i.e., $2^{\sum_i (n_m(i) + n_r(i))}$ combinations) as in Algorithm 1. Instead, we modify our objective in the dynamic slot configuration as there is no closed-form expression of the makespan.

Our goal now is, for the two concurrently running jobs (one in map phase and the other in reduce phase), to minimize the completion time of these two phases. Our intuition is to eliminate the two undesired cases mentioned above by aligning the completion of j_i 's map phase and j_{i-1} 's reduce phase, see Fig. 4(c). Briefly, we use the slot assignment as a tunable knob to change the level of parallelism of map or reduce tasks. When we assign more map slots, map tasks obtain more system resources and could be finished faster, and vice versa for reduce tasks. In the rest of this section, we first present our basic solution with the assumption of prior knowledge of job workload. Then, we describe how to estimate the workload in practice when it is not available. In addition, we present a feedback control-based solution to provide more accurate

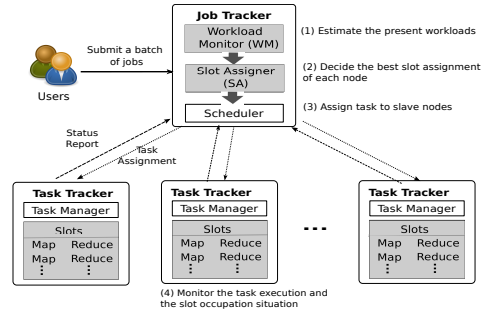


Fig. 5. The architecture overview of our design. The shade rectangles indicate our new/modified components in Hadoop.

estimation of the workload. Finally, we discuss the design of task scheduler in compliance with our solution.

A. Basic Sketch With Prior Knowledge of Workload

Assume the workload information is available, at the end of a task, Hadoop can obtain the value of the remaining workload for both map and reduce phases. Intuitively, we should assign more slots (resources) to the task type that has heavier remaining workload. Assume j_i and j_{i-1} are two active jobs and j_{i-1} is in reduce phase while j_i is in map phase. At the end of a task, we can get the number of remaining map tasks of j_i and remaining reduce tasks of j_{i-1} , indicated by $n'_m(i)$ and $n'_r(i-1)$. Let $w'_m(i)$ and $w'_r(i-1)$ represent the remaining workload of j_i 's map phase and j_{i-1} 's reduce phase, we have:

$$w'_m(i) = n'_m(i) \cdot \bar{t}_m(i), \quad w'_r(i-1) = n'_r(i-1) \cdot \bar{t}_r(i-1) \quad (2)$$

To align the completions of these two phases, the best parameters should satisfy the following condition:

$$\frac{n'_m(i)}{s_m} \cdot \bar{t}_m(i) = \frac{n'_r(i-1)}{s_r} \cdot \bar{t}_r(i-1) \Rightarrow \frac{w'_m(i)}{s_m} = \frac{w'_r(i-1)}{s_r} \quad (3)$$

Therefore, the number of map and reduce slots should be proportional to their remaining workloads as shown in Eq. 4-5,

$$s_m = \lfloor \frac{w'_m(i)}{w'_m(i) + w'_r(i-1)} \cdot S \rfloor, \quad (4)$$

$$s_r = S - s_m, \quad (5)$$

where s_m and s_r represent the target numbers of map and reduce slots respectively, and S is the total number of slots in the cluster which is configured based on system capacity. Furthermore, we introduce the upper bound s_m^h and the lower bound s_m^l for the map slots assignment. When the estimated value of s_m exceeds the bounds, we use the bound value as the new s_m value instead. In our design, s_m^l is set to be the number of nodes in the cluster (k) such that there is at least one map slot on each node at any time. Similarly, s_m^h is set to be equal to $S - s_m^l$ such that the reduce slots number in each node is always greater than or equal to 1. When a map or reduce task is finished, one slot becomes available. The Hadoop system calculates the values of s_m and s_r according to Eq. 4-5. If the current map slots are fewer than s_m , then the available slot will become a map slot and serve a map task. Otherwise, it turns to a reduce slot. With this setting, the current map and reduce phases could finish at approximately the same time with a high system resource utilization.

B. Workload Estimation

Our solution proposed above depends on the assumption of prior knowledge of workload information. In practice, workload can be derived from job profiles, training phase, or other empirical settings. In some applications, however, workload information may not be available or accurate. In this subsection, we present a method that estimates the workload during the job execution without any prior knowledge.

We use w'_m and w'_r to represent the remaining workload of a map or reduce phase, i.e., the summation of execution time of the unfinished map or reduce tasks. Note that we only track the map/reduce workloads of running jobs, but not the jobs waiting in the queue. Basically, the workload is calculated as the multiplication of the number of remaining tasks and the average task execution time of a job. Specifically, when a map or reduce task is finished, the current workload information needs to be updated, as shown in Algorithm 2, where $n'_m(i)/n'_r(i)$ is the number of unfinished map/reduce tasks of job j_i , and $\bar{t}_m(i)/\bar{t}_r(i)$ means the average execution time of finished map/reduce tasks from j_i . Note that the execution time of each finished task is already collected and reported to the JobTracker in current Hadoop systems. In addition, we use the Welford's one pass algorithm to calculate the average of task execution times, which incurs very low overheads on both time and memory space.

Algorithm 2 Workload Information Collector

```

if a map task of job  $j_i$  is finished then
  update the average execution time of a map task  $\bar{t}_m(i)$ 
   $w'_m(i) = \bar{t}_m(i) \cdot n'_m(i)$ 
if a reduce task of job  $j_i$  is finished then
  update the average execution time of a reduce task  $\bar{t}_r(i)$ 
   $w'_r(i) = \bar{t}_r(i) \cdot n'_r(i)$ 

```

C. Feedback Control-based Workload Estimation

In this subsection, we present an enhanced workload estimation algorithm to achieve more accurate workload information. Our previous analysis adopts an assumption that the execution time of a map or reduce task is similar, represented by the average values $\bar{t}_m(i)$ and $\bar{t}_r(i)$, respectively. They are also used for calculating the workload w_m and w_r . This estimation works well in systems where the slots assignment is fixed. In our system design, however, the slots assignment is dynamically changed, which affects the per task execution time in practice. Assigning more slots to one type of tasks may cause the contention on a particular system resource and lead to an increased execution time of each following task in the same type. For example, in ‘‘Case 2’’ described in Section II, when we use 1 map slot on each node, the average execution time of a map task is 18.5 sec. When we increase the number of map slots per node to 2, the average execution time of a map task becomes 23.1 sec with a 25% increase.

To overcome this issue, we have designed a feedback control based mechanism to tune the slots assignment. Under this mechanism, the slots assignment, s_m and s_r , is first calculated through Eq. 4-5. An additional routine is introduced

to periodically update the workload information based on newly captured average task execution times. If the workloads have changed, then the slots assignment will also be updated according to Eq. 6-7.

$$s_m = s_m + \lfloor \alpha \cdot \left(\frac{w'_m}{w'_m + w'_r} - \frac{w_m}{w_m + w_r} \right) \cdot S \rfloor, \quad (6)$$

$$s_r = S - s_m. \quad (7)$$

When the new estimated workloads, i.e., w'_m and w'_r , differ from the previous estimation, an integral gain parameter α is used to control the new assignment of slots based on the new estimation. The Hadoop system will iteratively calculate s_m and s_r (Eq. 6-7) until there is no change on these two parameters. The value of α is set to be 1.2 in our system such that the slots assignment could converge quickly.

D. Slot Assigner

The task assignment in Hadoop works in a heartbeat fashion: the TaskTrackers report slots occupation situation to the JobTracker with heartbeat messages; and the JobTracker selects tasks from the queue and assigns them to free slots. There are two new problems need to be addressed when assigning tasks under TuMM. First, slots of each type should be evenly distributed across the slave nodes. For example, when we have a new slot assignment $s_m = 5, s_r = 7$ in a cluster with 2 slave nodes, a 2:3/4:3 map/reduce slots distribution is better than the 1:4/5:2 map/reduce slots distribution in case of resource contention. Second, the currently running tasks may stick with their slots and therefore the new slot assignments may not be able to apply immediately. To address these problems, our slot assignment module (SA) takes both the slots assignment calculated through Eq. 6-7 and the situation of currently running tasks into consideration when assigning tasks.

The process of SA is shown in Algorithm 3. The SA first calculates the map and reduce slot assignments of slave node x (line 1), indicated by $s_m(x)$ and $s_r(x)$, based on the current values of s_m and s_r and the number of running tasks in cluster. Because of the flooring operation in line 1, the assigned slots ($s_m(x) + s_r(x)$) on node x may be fewer than the available slots (S/k). In lines 3–6, we increase either $s_m(x)$ or $s_r(x)$ to compensate slot assignment. Our decision is based on the deficit of current map and reduce slots (line 3), where s_m/s_r represent our target assignment and rt_m/rt_r are the number of current running map/reduce tasks. Eventually, we assign a task to the available slot in lines 7–10. Similarly, the decision is made by comparing the deficit of map and reduce tasks on node x , where $s_m(x)/s_r(x)$ are our target assignment and $rt_m(x)/rt_r(x)$ are the numbers of running tasks.

V. EVALUATION

A. Experimental Setup and Workloads

1) *Implementation*: We implemented our new scheme on the top of Hadoop Version 0.20.2. First, we added two new modules into the *JobTracker*: the Workload Monitor (WM) that is responsible to collect past workload information such as execution times of completed tasks and to estimate the

Algorithm 3 Slot Assigner

```
0: Input: Number of slave nodes in cluster:  $k$   
   Total numbers of running map/reduce tasks:  $rt_m, rt_r$ ;  
0: When receive heartbeat message from node  $x$  with the number  
   of running map/reduce tasks on node  $x$ :  $rt_m(x), rt_r(x)$ ;  
1: Initialize assignment of slots for node  $x$ :  
    $s_m(x) \leftarrow \lfloor s_m/k \rfloor, s_r(x) \leftarrow \lfloor s_r/k \rfloor$ ;  
2: if  $(s_m(x) + s_r(x)) < S/k$  then  
3:   if  $(s_m - rt_m) > (s_r - rt_r)$  then  
4:      $s_m(x) \leftarrow s_m(x) + 1$ ;  
5:   else  
6:      $s_r(x) \leftarrow s_r(x) + 1$ ;  
7:   if  $(s_m(x) - rt_m(x)) > (s_r(x) - rt_r(x))$  then  
8:     assign a map task to node  $x$ ;  
9:   else  
10:    assign a reduce task to node  $x$ .
```

workloads of currently running map and reduce tasks and the Slot Assigner (SA) which uses the estimated information received from WM to adjust the slot ratio between map and reduce for each slave node. The *JobTracker* with these additional modules will then assign tasks to a slave node based on the adjusted slot ratio and the current slot status at that particular node. In addition, we modified the *TaskTracker* as well as the *JvmManager* at each slave node to check the number of individual map and reduce tasks running on that node based on the new slot ratio received from the *JobTracker*. The architecture overview of this new Hadoop framework is shown in Fig. 5.

2) *Workloads*: We choose five representative Hadoop benchmarks from Purdue MapReduce Benchmarks Suite [12]:

- *Inverted Index*: take text documents as input and generate word to document indexing.
- *Histogram Rating*: take the movie rating data as input and calculate a histogram of input data.
- *Word Count*: take text documents as input and count the occurrence of each word.
- *Classification*: take the movie rating data as input and classify the movies into one of the predefined clusters.
- *Grep*: take text documents as input and search for a pattern in the files.

In addition, we use a 10GB movie rating data [12] that consists of user ranking information and a 7GB wiki category links data [13] that includes the information about wiki page categories, as the input to the above five benchmarks.

3) *Hadoop Cluster*: All the experiments are conducted in a Hadoop cluster which consists of 5 *m1.xlarge* Amazon EC2 instances. Specifically, we have one master node and four slave nodes in the cluster. The number of slots which can be available on each slave node is set as 4 since an *m1.xlarge* instance at Amazon EC2 has 4 virtual cores.

B. Performance Evaluation

In this section, we evaluate the performance of TuMM in terms of the makespan of a batch of MapReduce jobs. We first consider the simple workloads which consist of jobs from a single MapReduce benchmark and then validate the robustness

of our approach with a mixed workload that is a combination of different MapReduce benchmarks.

1) *Simple Workloads*: We here conduct a set of experiments such that in each experiment 3 Hadoop jobs from one of the above benchmarks (see Section V-A2) are waiting for service. We remark that such a simple workload is often found in real systems as the same Hadoop jobs may be executed repeatedly to process similar or different input data sets. In our experiments, three Hadoop jobs use the same data set as the input. Furthermore, as the comparisons, we evaluate the performance under the static slot ratios for map and reduce. Since all the slave nodes normally have the same slot ratio in current Hadoop implementations, With our setting in the evaluation (i.e., total number of slots per node is 4), there are three static configuration alternatives, i.e., 1:3, 2:2 and 3:1, for a Hadoop cluster. So we enumerate all these three possible settings for the comparison with our solution.

Fig. 6 shows the makespans (i.e., the completion lengths) of a given set when we have different slot configurations. We first observe that the performance varies a lot under three static slot settings. For example, the *Inverted Index* jobs experience the fastest makespan when the slot ratio is equal to 1:3. In contrast, the *Histogram Rating* jobs achieve better performance when we assign more slots to their map tasks, e.g., with slot ratio of 3:1. We also observe that TuMM always yields the best performance, i.e., the shortest makespan, for all the five Hadoop benchmarks. We interpret this effect as the result of dynamic slot ratio adjustments enabled by TuMM.

Compared to the slot ratio of 2:2, our approach in average achieves about 20% relative improvement in the makespan. Moreover, such improvement becomes more visible when the workloads of map and reduce tasks become more unbalanced. For example, the makespan of the *Inverted Index* jobs is reduced by 28% where these jobs have their reduce phases longer than their map phases.

2) *Mixed Workloads*: In the previous experiments, each workload only contains jobs from the same benchmark. Now, we consider a more complex workload, which mixes jobs from different Hadoop benchmarks. Reducing the makespan for such a mixed workload thus becomes non-trivial. One solution to tackle this problem is to shuffle the execution order of these jobs. For example, the classic Johnson's algorithm [14] that was proposed for building an optimal two-stage job schedule, could be applied to process a set of Hadoop jobs and minimize the makespan of a given set as well. However, this algorithm needs to assume a priori knowledge of the exact execution times of each job's map and reduce phases, which unfortunately limits the adoption of this algorithm in real Hadoop systems. Moreover, for some cases, it may not be feasible to change the execution order of jobs, especially when there exists dependency among jobs or some of them have high priority to be processed first.

To address the above issues, our solution leverages the knowledge of the completed tasks to estimate the execution times of the currently running tasks and reduces the makespan of a set of jobs by dynamically adjusting the slot assignments

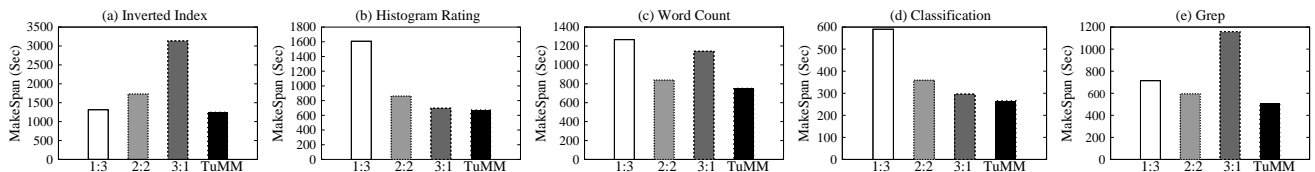


Fig. 6. Makespans of five Hadoop applications under TuMM and three static slot configurations.

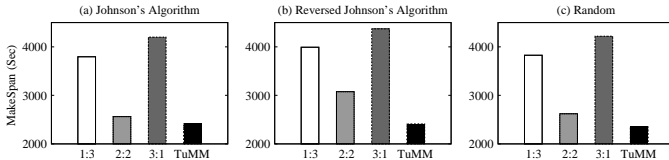


Fig. 7. Makespans of a mixed workload under TuMM and three static slot configurations. Three execution orders are also considered: (a) a sequence follows Johnson’s algorithm, (b) a sequence with reversed order of Johnson’s algorithm, and (c) a random sequence.

for map and reduce tasks. As a result, TuMM does not need to change the execution order of jobs and does not need to know the exact task execution times in advance, either.

We generate the mixed workload for our experiments by randomly choosing 10 jobs from 5 different Hadoop benchmarks. In order to investigate the impact of job execution order, we also consider three different execution sequences, including (1) a sequence generated by Johnson’s algorithm which can be considered as the optimal case in terms of the makespan; (2) a sequence that is inverse to the first one and can be considered as the worst case; and (3) a sequence that is random. Similarly, we evaluate the performance (i.e., makespan) under TuMM and three static slot configurations.

Fig. 7 shows the makespans of the 10 jobs in the mixed workload. We first observe that among three static settings, the slot ratio of 2:2 always achieves the best performance under three different execution orders. This is because the overall workloads of map tasks and reduce tasks from the 10 jobs are well balanced. We also notice that with a fixed number of slots per node, different job execution orders could yield different makespans. While our solution always achieves the best performance and the impact of execution sequence on our solution’s performance becomes less visible. This means that no matter what the execution order is, TuMM can always serve the jobs with the shortest makespans. That is, our approach allows to improve the performance in terms of makespan without changing the execution order of jobs.

To better understand how TuMM uses the slot ratio as a tunable knob to improve the makespan, we further plot the task execution times for each job as well as the transient slot assignments in Fig. 8, where the plots in the first row depict the running period of each task from the 10 jobs while the plots in the second row illustrate how the slot assignments change across time. As shown in Fig. 8, TuMM dynamically adjusts the slot assignments to map and reduce tasks based on the estimated workload information. For example, in the first 1200 seconds of Fig. 8-(2), TuMM attempts to assign more slots to reduce tasks. Then, in the later 1200 seconds, TuMM turns to allow more available map slots on each node. This is

because the Johnson’s algorithm shuffles the order of 10 jobs such that all the reduce intensive jobs such as *Inverted Index* and *Grep* run before the map intensive jobs, e.g., *Histogram Rating* and *Classification*. The only exception is the first 100s where most of the slots are assigned to map tasks even though the running job actually has reduce intensive workloads. That is because TuMM does not consider the reduce workloads of this job in the first 100 seconds until its map tasks are finished. Fig. 8-(1) shows the corresponding task execution times under TuMM. It is obvious that each job’s reduce phase successfully overlaps with the map phase of the following job and the makespan of 10 jobs is then shortened compared to the static settings.

In summary, TuMM achieves non-negligible improvements in makespan under both simple workloads and mixed workloads. By leveraging the history information, our solution accurately captures the changes in map and reduce workloads and adapts to such changes by adjusting the slot assignments for these two types of tasks. Furthermore, different job execution orders do not affect TuMM’s performance. That is, our solution can still reduce the makespan without changing the execution order of a given set of jobs.

VI. RELATED WORKS

An important direction for improving the performance of a Hadoop system is job scheduling. The default FIFO scheduler does not work well in a shared cluster with multiple users and a variety of jobs. Fair [15] and Capacity [16] schedulers were proposed to ensure that each job can get a proper share of the available resources; and Quincy [5] addressed the scheduling problem with locality and fairness constraints. Recently, Zaharia et al. [3] proposed a delay scheduling to further improve the performance of the Fair scheduler by increasing data locality. Verma et al. [4] introduced a heuristic to minimize the makespan of a set of independent MapReduce jobs by applying the classic Johnson’s algorithm.

Another category of schedulers further consider user-level goals while improving the performance. ARIA, a deadline aware scheduler, was recently proposed in [6], which always schedules a job with the earliest deadline and uses the Lagrange’s method to find the minimum number of slots for each job in order to meet the predefined deadline. Similarly, Polo et al. [7] estimated the task execution times based on the average execution times of the completed tasks instead of the job profiles. Task execution times were then used to calculate the number of slots that a job needed to meet its deadline. Although these deadline aware schedulers support user-level goals, their techniques are still based on static slot

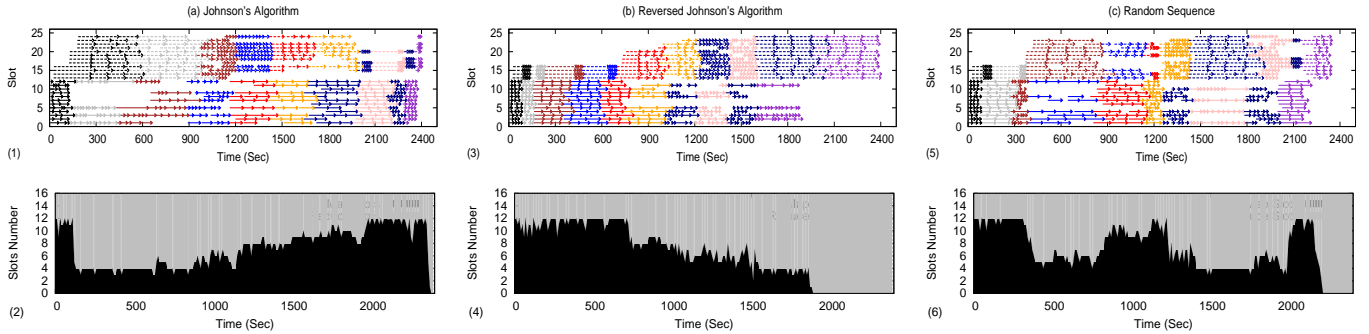


Fig. 8. Illustrating task execution times and slot assignments across time under TuMM, where the job execution sequence is (a) generated by Johnson's algorithm; (b) inverse to the first one; and (c) random. In the plots at the second row, black (resp. gray) areas represent the number of available map (resp. reduce) slots in the cluster.

configurations, i.e., having a fixed number of map slots and reduce slots per node throughout the lifetime of a cluster.

Finally, resource aware management is another important direction for improving performance in Hadoop. RAS [10] leverages existing profiling information to dynamically determine the number of job slots and their placement in the cluster. The goal of this approach is to maximize the resource utilization of the cluster and to meet job completion time deadlines. More recently, [11] introduces a local resource manager at each TaskTracker to detect task resource utilization and predict task finish time, and a global resource manager at the JobTracker to coordinate the resource assignments to each task; and [9] addresses the cluster resource utilization problem by developing a dynamic split model of resource utilization. The Hadoop community recently released Next Generation MapReduce (NGM) [8], the latest architecture of Hadoop MapReduce, which replaces the fixed-size slot with a resource container that works in a fine-grained resource level. In this work, we rely on the history task execution times instead of the profiling of fine-grained resource usage to dynamically adjust the number of map and reduce slots at each node. The main objective of our work is to reduce the completion length (i.e., makespan) of a set of MapReduce jobs. Our work is complementary to the above techniques.

VII. CONCLUSION

In this paper, we presented a novel slot management scheme, named TuMM, to enable dynamic slot configuration in Hadoop. The main objective of TuMM is to improve resource utilization and reduce the makespan of multiple jobs. To meet this goal, the presented scheme introduces two main components: *Workload Monitor* periodically tracks the execution information of recently completed tasks and estimates the present workloads of map and reduce tasks and *Slot Assigner* dynamically allocates the slots to map and reduce tasks by leveraging the estimated workload information. We implemented TuMM on the top of Hadoop v0.20.2 and evaluated this scheme by running representative MapReduce benchmarks in a cluster at Amazon EC2. The experimental results demonstrate up to 28% reduction in the makespans and 20% increase in resource utilizations. The effectiveness

and the robustness of TuMM are validated under both simple workloads and a more complex mixed workload. In the future, we plan to extend the current work to enable performance improvement in a Hadoop system with decentralized JobTrackers and evaluate TuMM in a large-scaled experimental test-bed.

REFERENCES

- [1] J. Dean, S. Ghemawat, and G. Inc, "Mapreduce: simplified data processing on large clusters," in *OSDI'04*, 2004.
- [2] Apache Hadoop. [Online]. Available: <http://hadoop.apache.org/>
- [3] M. Zaharia, D. Borthakur, J. S. Sarma *et al.*, "Delay scheduling: A simple technique for achieving locality and fairness in cluster scheduling," in *EuroSys'10*, 2010.
- [4] A. Verma, L. Cherkasova, and R. H. Campbell, "Two sides of a coin: Optimizing the schedule of mapreduce jobs to minimize their makespan and improve cluster performance," in *MASCOTS'12*, Aug 2012.
- [5] M. Isard, Vijayan Prabhakaran, J. Currey *et al.*, "Quincy: fair scheduling for distributed computing clusters," in *SOSP'09*, 2009, pp. 261–276.
- [6] A. Verma, Ludmila Cherkasova, and R. H. Campbell, "Aria: Automatic resource inference and allocation for mapreduce environments," in *ICAC'11*, 2011, pp. 235–244.
- [7] J. Polo, D. Carrera, Y. Becerra *et al.*, "Performance-driven task co-scheduling for mapreduce environments," in *NOMS'10*, 2010.
- [8] Next generation mapreduce scheduler. [Online]. Available: <http://goo.gl/GACMM>
- [9] X. W. Wang, J. Zhang, H. M. Liao, and L. Zha, "Dynamic split model of resource utilization in mapreduce," ser. DataCloud-SC '11, 2011.
- [10] J. Polo, C. Castillo, D. Carrera *et al.*, "Resource-aware adaptive scheduling for mapreduce clusters," in *Proceedings of the 12th ACM/IFIP/USENIX international conference on Middleware*, 2011.
- [11] B. Sharma, R. Prabhakar, S.-H. Lim *et al.*, "Mrorchestrator: A fine-grained resource orchestration framework for mapreduce clusters," in *CLOUD'12*, 2012.
- [12] Purdue mapreduce benchmarks suite. [Online]. Available: <http://web.ics.purdue.edu/~fahmad/benchmarks.htm>
- [13] Wiki data sets. [Online]. Available: <http://dumps.wikimedia.org/>
- [14] S. M. Johnson, "Optimal two- and three-stage production schedules with setup times included," *Naval Research Logistics Quarterly*, vol. 1, no. 1, pp. 61–68, 1954.
- [15] M. Zaharia, D. Borthakur, J. S. Sarma *et al.*, "Job scheduling for multi-user mapreduce clusters," University of California, Berkeley, Tech. Rep., Apr. 2009.
- [16] Capacity scheduler. [Online]. Available: http://hadoop.apache.org/common/docs/r1.0.0/capacity_scheduler.html