

SCache: Efficient and General Shuffle Management for DAG Computing Frameworks

Rui Ren, Chunghsuan Wu, Zhouwang Fu, Tao Song, Zhengwei Qi, *Member, IEEE*, Haibing Guan, *Member, IEEE*,

Abstract—In large-scale data-parallel analytics, *shuffle*, or the cross-network read and aggregation of partitioned data between tasks with data dependencies, usually brings in large network transfer overhead. Due to the dependency constraints, execution of those descendant tasks could be delayed by logy shuffles. To reduce shuffle overhead, we present *SCache*, an open source plug-in system that particularly focuses on shuffle optimization in frameworks defining jobs as *directed acyclic graphs* (DAGs). During shuffle optimization process, we propose a performance model called *Framework Resources Quantification* (FRQ) model. FRQ model displays utilization of resources in time dimension and calculate execution time of computing jobs on different DAG frameworks. We use FRQ model to assist in analyzing shuffle process and verify *SCache* shuffle optimization by mathematics. By extracting and analyzing the DAGs and shuffle dependencies prior to the actual task execution, *SCache* can take full advantage of the system memory to accelerate the shuffle process. Meanwhile, it adopts heuristic-MinHeap scheduling combining with shuffle size prediction to pre-fetch shuffle data and balance the total size of data that will be processed by each descendant task on each node. We have implemented *SCache* as the external shuffle service and co-scheduler on both Apache Spark and Apache Hadoop. The performance of *SCache* is evaluated with both simulations and testbed experiments on a 50-node Amazon EC2 cluster. Those evaluations have demonstrated that, by incorporating *SCache*, the shuffle overhead of Spark can be reduced by nearly 89%, and the overall completion time of TPC-DS queries improves 40% on average. On Apache Hadoop, *SCache* optimize end-to-end Terasort completion time by 15%.

Index Terms—Distributed DAG frameworks, Shuffle, Optimization

1 INTRODUCTION

RECENT years have witnessed widespread use of sophisticated frameworks such as Dryad [1], Spark [2], and Apache Tez [3]. Despite the differences among data-intensive frameworks, their communication is always structured as a shuffle phase, which takes place between successive computation stages. Such shuffle phase places significant burden for both the disk and network I/O, thus heavily affecting the end-to-end application performance. For instance, a MapReduce trace analysis from Facebook shows that shuffle accounts for 33% of the job completion time on average, and up to 70% in shuffle-heavy jobs [4].

Although continued efforts of performance optimization have been made among a variety of computing frameworks [5], [6], [7], [8], [9], [10], the shuffle is often poorly optimized in practice. In particular, we observe that one major deficiency lies in a lack of fine-grained, coordinated management among different system resources. As Figure 1 shows, the *shuffle write* is responsible for writing intermediate results to disk, which is attached to the tasks in ancestor stage (i.e., map task). And the *shuffle read* fetches intermediate results from *remote* disks through network, which is commonly integrated as part of the tasks in descendant stage (i.e., reduce task). Once scheduled, a fixed bundle of resources (i.e., CPU, memory, disk and network) named *slot* is assigned to each of the computation task, and

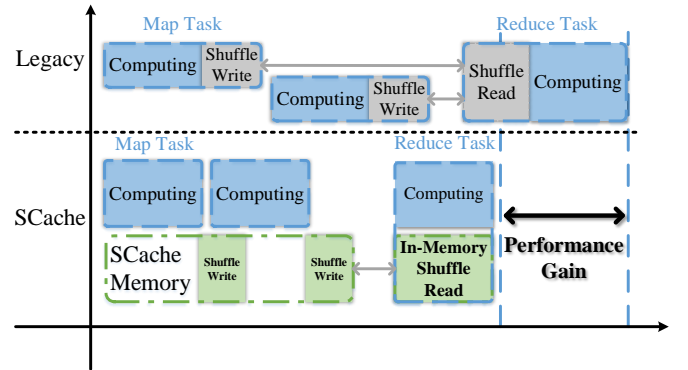


Fig. 1: Workflow Comparison between Legacy DAG Computing Frameworks and Frameworks with *SCache*

the resources are released only after the task finishes. Such task aggregation together with the coarse-grained scheduling effectively simplifies task management. However, since each phase receives a fixed bundle of resources, attaching the *I/O intensive* shuffle phase to the *CPU/memory intensive* computation phase results in a poor multiplexing between computational and I/O resources.

Moreover, the shuffle read phase introduces all-to-all communication pattern across the network, and such network I/O procedure is also poorly coordinated. Note that the shuffle read phase starts fetching data only after all the data from its ancestor stage by default. As a result, all the corresponding reduce tasks start fetching shuffle data almost simultaneously. Such synchronized network communication causes a burst demand for network I/O,

- R. Ren ...E-mail: renrui@sjtu.edu.cn
- W. Fu and C. Wu ...E-mail: {wuchunghsuan, frankfzw}@sjtu.edu.cn
- T. Song ...E-mail: songt333@sjtu.edu.cn
- Z. Qi ...E-mail: qizhwei@sjtu.edu.cn
- H. Guan is with Shanghai Key Laboratory of Scalable Computing and Systems, and Department of Computer Science, Shanghai Jiao Tong University, Shanghai, China. E-mail: hbguan@sjtu.edu.cn

which in turn greatly enlarges the shuffle read completion time. To desynchronize the network communication, an intuitive way is to launch some tasks in the descendent stage earlier such as *early start* from Apache Hadoop¹. However, such early start is by no means a panacea. This is because although the reduce tasks can start fetching data earlier, the computation can only take place after all the data is ready. Since each phase receives a fixed bundle of resources, starting a reduce task early always introduces an unnecessary early allocation of slot.

To make things worse, we note that the above deficiencies generally exist in most of the DAG computing frameworks. As a result, even we can effectively resolve the above deficiencies by modifying one framework, updating one application at a time is impractical given the sheer number of computing frameworks available today.

In order to visually analyze above-mentioned resources scheduling, we propose a performance model called *Framework Resources Quantification*(FRQ) model. FRQ model quantifies computing and I/O resources and displays resources scheduling strategy of DAG framework in time dimension. We use FRQ model to assist in analyze the deficiencies of resources scheduling and optimize it. FRQ model has five parameters: *Input Data Size*, *Data Conversion Rate*, *Computation Round Number*, *Computation Speed*, and *Shuffle Speed*. According to the above five parameters and the scheduling strategy of the DAG framework, FRQ model is able to calculate the execution time of each phase of a computing job. Take Apache Hadoop MapReduce as a simple DAG computing example, we use FRQ to model each phases of the computing, including Map, Shuffle, and Reduce. By revealing the relationships between the various phases, FRQ model assist us in discovering the irrationality of resource scheduling during computing. I/O resource in Map phase is not fully utilized. Furthermore, due to coupling with reduce phase, shuffle phase delay the execution time of reduce phase. This situation will be amplified in shuffle-sensitive tasks.

Can we efficiently optimize the data shuffling without significantly changing every framework? In this paper, we answer this question in the affirmative with S(huffle)Cache, an open source plug-in system which provides a shuffle-specific optimization for different DAG computing frameworks.

Specifically, SCache takes over the whole shuffle phase from the underlying framework by providing a cross-framework API for both shuffle write and read. SCache's effectiveness lies in the following two key ideas. First, SCache decouples the shuffle read and write from both map and reduce tasks. Such decoupling effectively enables fine-grained resource management and better multiplexing between the computational and I/O resources. In addition, SCache pre-schedules the reduce tasks and pre-fetches the shuffle data to the location of the reduce tasks without launching them. Such pre-scheduling and pre-fetching effectively desynchronize the network I/O operation, while avoiding the waste of computational resources compared to the early start mechanism.

The workflow of DAG framework with SCache is presented in Figure 1. SCache hijacks the intermediate data of a map task in memory space. The disk operation is skipped and the slot is released after memory copy. The in-memory intermediate data is immediately pre-fetched through network after pre-scheduling. By releasing the slot earlier and starting the network transfer ahead of reduce tasks, SCache can help the DAG framework achieve a significant performance gain. A by-product optimization of pre-scheduling is that SCache can provide a more balanced load for each node and further benefit the reduce stage by avoiding data skew.

The main challenge to achieve this optimization is *pre-scheduling reduce tasks*. It is not critical for the simple DAG computing such as Hadoop MapReduce [11]. Unfortunately the complexity of DAG can amplify the defects of naïve pre-scheduling schemes. In particular, randomly assign reduce tasks might result in a collision of two heavy tasks on one node. This collision can aggravate data skew, thus hurting the performance. To address this challenge, we propose a heuristic scheme to predict the shuffle output distribution and pre-schedule reduce tasks (Section 3).

The second challenge is the *in-memory data management*. To prevent shuffle data touching the disk, SCache leverages extra memory to store the shuffle data. However, the memory is a precious resource for DAG computing. To minimize the reserved memory while maximizing the performance gain of optimization and memory utilization, we propose two constraints: all-or-nothing and context-aware (Section 4.2).

We have implemented SCache, a customized Apache Spark [12] and a customized Apache Hadoop. We have also designed a performance model called *Framework Resources Quantification*(FRQ) model. We use FRQ model to analyze the shuffle optimization of SCache and calculate the execution time of each phase of computing job. The performance of SCache is evaluated with both simulations and testbed experiments on a 50-node Amazon EC2 cluster on both Apache Spark and Apache Hadoop. On Apache Spark, we conduct basic test GroupByTest. We also evaluate Terasort² benchmark and standard workloads like TPC-DS³ for multi-tenant modeling. On Apache Hadoop, we focus on Terasort benchmark. In a nutshell, SCache can eliminate explicit shuffle process by at most 89% in varied application scenarios. More impressively, SCache reduces 40% of overall completion time of TPC-DS queries on average on Apache Spark. On Apache Hadoop, SCache optimize end-to-end Terasort completion time by 15%.

2 BACKGROUND AND OBSERVATIONS

In this section, we first study the typical shuffle characteristics (2.1), and then spot the opportunities to achieve shuffle optimization (2.2).

2.1 Characteristic of Shuffle

In large scale data parallel computing, shuffle is designed to achieve an all-to-all data blocks transfer among nodes. For

1. <http://hadoop.apache.com/>

2. <https://github.com/ehiggs/spark-terasort>

3. <http://www.tpc.org/tpcds/>



Fig. 2: CPU utilization and I/O throughput of a node during a Spark single shuffle application

a clear illustration, we use *map tasks* to define the tasks that generate shuffle data and use *reduce tasks* to define the tasks that consume shuffle data.

Overview of shuffle process. Each map task partitions the result data (key, value pair) into several buckets according to the partition function (e.g., hash). The total number of buckets equals the number of tasks in the next step. The shuffle process can be further split into two parts: *shuffle write* and *shuffle read*. The partitioned shuffle output data will be spilled to local persistent storage during shuffle write. Shuffle read starts at the beginning of reduce tasks. It fetches data as reduce input from both remote and local storage.

Impact of shuffle process. Shuffle process is I/O intensive, which might introduces a significant latency to the application. Reports show that 60% of MapReduce jobs at Yahoo! and 20% at Facebook are shuffle intensive workloads [13]. For those shuffle intensive jobs, the shuffle latency may even dominate Job Completion Time (JCT). For instance, a MapReduce trace analysis from Facebook shows that shuffle accounts for 33% JCT on average, up to 70% in shuffle intensive jobs [4].

2.2 Observations

Of course, shuffle is unavoidable in a DAG computing process. But can we mitigate or even remove the overhead of shuffle? To find the answers, we run some typical Spark applications in a 5-node EC2 cluster with m4.xlarge. We then measure the CPU utilization, I/O throughput, and tasks execution information of each node. Here we present the trace of one node running Spark *GroupByTest* in Figure 2 as an example. This job has 2 rounds of tasks for each node. We have marked the *execution* phase as from the launch time of the first task to the execution finish timestamp of the last one. The *shuffle write* phase is marked from the timestamp of the beginning of the first partitioned data write. The *shuffle read and execution* phase is marked from the start of the first reduce launch timestamp.

2.2.1 Coarse Granularity Resource Allocation

In Figure 2, when a slot is scheduled to a task, it will not be released until the the task completes *shuffle write*. On the reduce side, the network transfer of shuffle data introduces an explicit I/O delay during *shuffle read*. Meanwhile, shuffle is I/O intensive. Both *shuffle write* and *shuffle read* occupy the slot without involving much CPU. The current coarse slot — task mapping results in an inconsistency between resource demands and slot allocation thus decreases the resource

utilization. To break this inconsistency, a finer granularity resource allocation scheme must be provided.

2.2.2 Synchronized Shuffle Read

Almost all reduce tasks start *shuffle read* simultaneously. The synchronized *shuffle read* requests cause a burst of network traffic. As shown in Figure 2, the data transfer put a great deal of stress on network bandwidth, which may results in network congestion and further hurts the performance of reduce stage.

2.2.3 Inefficient Persistent Storage Operation

At first, both shuffle write and read are tightly coupled with task execution, which results in a blocking style I/O operation. This blocking I/O operation along with synchronized shuffle read may introduce significant latency, especially in a I/O performance bounded cluster. Besides, the legacy of spilling shuffle data to disk is too conservative in modern cluster with large memory. Compared to input dataset, the size of shuffle data is relatively small. For example, shuffle size of Spark Terasort is less than 25% of input data. The data reported in [14] also shows that the amount of data shuffled is less than input data, by as much as 10%-20%. On the other hand, numbers of memory based distributed storage system have been proposed [7], [15], [16] to move data back to memory. We argue that the memory capacity is large enough to store the short-living shuffle data with cautious management.

2.2.4 Multi-round Tasks Execution

Both experience and DAG framework manuals recommend that multi-round execution of each stage will benefit the performance of applications. For example, Hadoop MapReduce Tutorial¹ suggests that $10-100$ maps per-node and 0.95 or $1.75 \times \text{no. of nodes} \times \text{no. of maximum container per-node}$ seem to be the right level of parallelism. Spark Configuration also recommends $2 - 3$ tasks per CPU core². Since the shuffle data becomes available as soon as the end of map task execution. At the same time, the network is idle during the map stage (network utilization during map stage in Figure 2). If the destination host of the shuffle data can be predicted in priori, the property of multi-round can be leveraged to overlap the *shuffle read* operation.

Based on these observations, it is straightforward to come up with an optimization that starts shuffle read ahead of reduce stage to overlap the I/O operations in multi-round DAG computing tasks, and uses memory to store the shuffle data for further decreasing shuffle overhead. To achieve this optimization:

- Shuffle process should be decoupled from task execution to achieve a fine granularity scheduling scheme.
- Reduce tasks should be pre-scheduled without launching to achieve shuffle data pre-fetching.
- Shuffle process should be taken over and managed outside DAG frameworks to achieve a cross-framework optimization.

1. <http://hadoop.apache.org/docs/current/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html/>
 2. <http://spark.apache.org/docs/1.6.2/configuration.html>

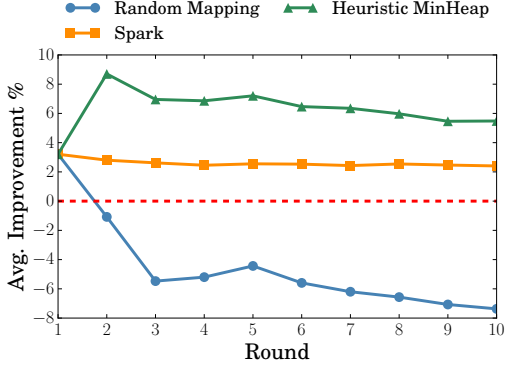


Fig. 3: Stage Completion Time Improvement of OpenCloud Trace

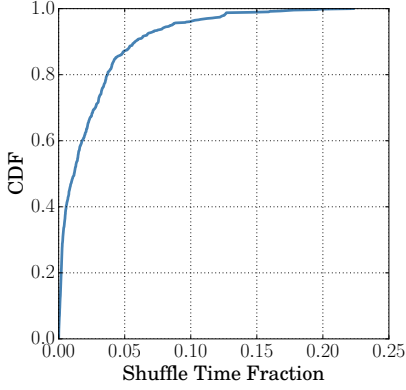


Fig. 4: Shuffle Time Fraction CDF of OpenCloud Trace

3 SHUFFLE OPTIMIZATION

This section presents the detailed methodologies to achieve three design goals. The out-of-framework shuffle data management is used to decouple shuffle from execution and provide a cross-framework optimization. Two heuristic algorithms (Algorithm 1, 2) are used to achieve shuffle data pre-fetching without launching tasks.

3.1 Decouple Shuffle from Execution

During map tasks, the partitioner takes a set of key-value pairs as input and calculates the partition number for each of them by applying pre-defined the partition function. After that, it stores all key-value pairs in the corresponding data blocks. Each of the block contains the key-value pairs for one reduce partition. At the same time, blocks will be spilled to disk. In order to avoid blocking the slot by disk operation, we use memory copy to hijack shuffle data from map tasks. By doing this, a slot can be released as soon as a task finishes CPU intensive computing. From the perspective of reduce task, shuffle read is decoupled by pre-fetching shuffle data to local SCache client before reduce tasks start.

3.2 Pre-schedule with Application Context

The pre-scheduling and pre-fetching start when the collected shuffle data exceeds the threshold. This is the most critic step toward the optimization. The task — node mapping is not determined until tasks are scheduled by the

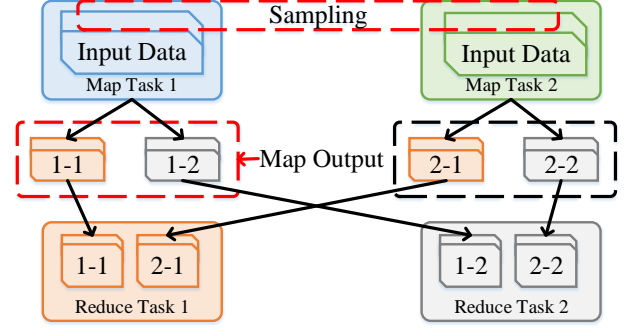


Fig. 5: Shuffle Data Prediction

scheduler of DAG framework. Once the tasks are scheduled, the slots will be occupied to launch them. On the other hand, the shuffle data cannot be pre-fetched without the readiness of task — node mapping. To get rid of this dilemma, we propose a co-scheduling scheme with two heuristic algorithms (Algorithm 1, 2). That is, the task — node mapping is established ahead of DAG framework scheduler, and it is enforced to DAG scheduler before the real scheduling.

3.2.1 Problem of Random Mapping

The simplest way of pre-scheduling is mapping tasks to different nodes randomly. It only guarantees that each node run same number of tasks. As shown in Figure 3, we use traces from OpenCloud¹ for the simulation to evaluate the impact of different pre-scheduling algorithms. The baseline (red dot line in Figure 3) is the stage completion time with Spark default scheduling algorithm. And then we remove the shuffle read time of each task, and run the simulation under three different schemes: random mapping, Spark FIFO, and our heuristic MinHeap. Note that most of the traces from OpenCloud is shuffle-light workload as shown in Figure 4. The average shuffle read time is 3.2% of total reduce completion time.

Random mapping works well when there is only one round of tasks. But the performance collapses as the round number grows. It is because that data skew commonly exists in data-parallel computing [17], [18], [19]. Several heavy tasks might be assigned on the same node. This collision then slows down the whole stage and makes the performance even worse than the baseline. In addition, randomly assigned tasks also ignore the data locality between shuffle map output and shuffle reduce input, which might introduces extra network traffic in cluster.

3.2.2 Shuffle Output Prediction

The problem of random mapping was obviously caused by application context (e.g., shuffle data size) unawareness. Note that the optimal schedule decision can be made under the awareness of shuffle dependencies number, partition number, and shuffle size for each partition. The first two of them can be easily extracted from DAG information. The scheduling can be made with the “prediction” of shuffle size.

According to the DAG computing process, the shuffle size of each reduce task is decided by input data, map

1. <http://ftp.pdl.cmu.edu/pub/datasets/hla/dataset.html>

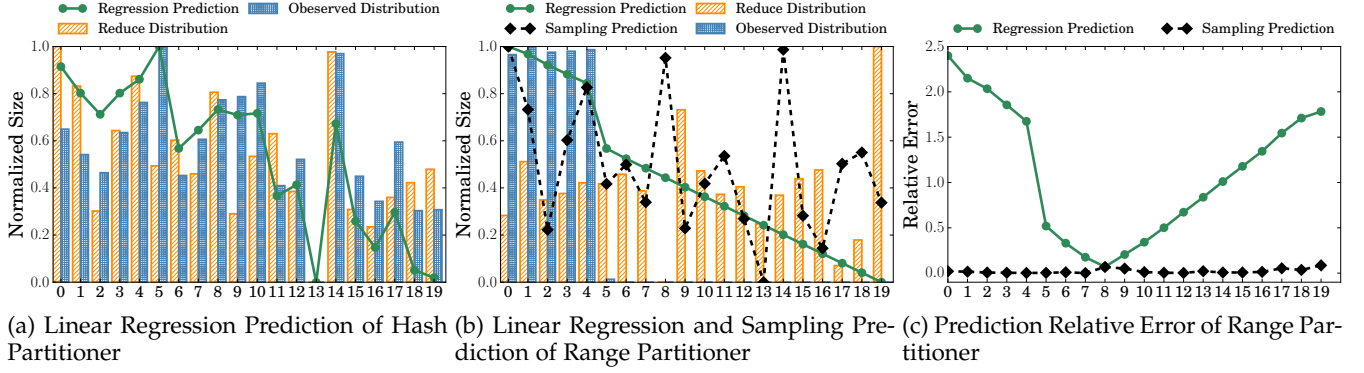


Fig. 6: Reduction Distribution Prediction

task computation, and hash partitioner. Each map task produces a data block for each reduce task. The size of each reduce partition can be calculated $reduceSize_i = \sum_{j=0}^m BlockSize_{ji}$ (m is the number of map tasks). $BlockSize_{ji}$ represents the size of block which is produced by map $task_j$ for reduce $task_i$ (e.g., block '1-1' in Figure 5).

For the simple DAG applications such as Hadoop MapReduce [11], the $BlockSize_{ji}$ can be predicted with decent accuracy by liner regression model based on observation that the ratio of map output size and input size are invariant given the same job configuration [20], [21].

But the sophisticated DAG computing frameworks like Spark introduce more uncertainties. For instance, the customized partitioner might bring large inconsistency between observed map output blocks distribution and the final shuffle data distribution. In Figure 6, we use different datasets with different partitioners, and normalize the distribution to 0 – 1 to fit in one figure. In Figure 6a, we use a random input dataset with the hash partitioner. In Figure 6b, we use a skew dataset with the range partitioner of Spark [12]. The observed map outputs are randomly picked. As we can see, in hash partitioner, the distribution of observed map output is close to the final reduce input distribution. The prediction results also turn out to be good. However, this inconsistency results in a deviation in linear regression model. To handle this inconsistency, we introduce another methodology named weighted reservoir sampling. The $BlockSize_{ji}$ can be calculated by

$$BlockSize_{ji} = InputSize_j \times \frac{sample_i}{s \times p} \quad (1)$$

$sample_i = \text{number of samples for } reduce_i$

For each partition of map task, we use classic reservoir sampling to randomly pick $s \times p$ of samples, where p is the number of reduce tasks and s is a tunable number. After that, the map function is called locally to process the sampled data (*sampling* in Figure 5). The final sampling outputs are collected with the $InputSize$ of each map partition which is used as the weight for each set of samples.

In Figure 6b, we set s equals 3, the result of sampling prediction is much better even in a very skew scenario. The variance of the normalization between sampling prediction and reduce distribution is because the standard deviation of the prediction result is relatively small compared to the average prediction size, which is 0.0015 in this example.

Figure 6c further proves that the sampling prediction can provide precise result even in the dimension of absolute shuffle partition size. On the opposite, the result of linear regression comes out with large relative error.

To avoid extra overhead, the sampling prediction will be triggered only when the range partitioner or customized non-hash partitioner occurs. We will show the detail evaluation of sampling in the Section 6.

During the phase of shuffle output prediction, the composition of each reduce partition is calculated as well. We define $prob_i$ as

$$prob_i = \max_{0 \leq j \leq m} \frac{BlockSize_{ji}}{reduceSize_i} \quad (2)$$

$m = \text{number of map tasks}$

This parameter is used to achieve a better locality while performing shuffle scheduling.

3.2.3 Heuristic MinHeap Scheduling

In order to balance load on each node while reducing the network traffic, we present a heuristic MinHeap scheduling algorithm (Algorithm 1) for single shuffle. Heuristic MinHeap can be divided into two rounds. In the first round (i.e., the first *while* in *SCHEDULE*), the reduce tasks are first sorted by size in a descending order. For hosts, we use a min-heap according to size of assigned tasks to maintain the priority. So that the tasks can be distributed evenly in the cluster. In the second round, the task — node mapping will be adjusted according to the locality. The *SWAP_TASKS* will only be triggered when the *host_id* of a reduce task is not equal the *assigned_id*. For a task which contains at most $prob$ data from *host*, the normalized probability *norm* is calculated as a bound of performance degradation. We set maximum *upper_bound* of performance degradation equals to 10% that can be traded for locality (in extreme skew scenarios). Inside the *SWAP_TASKS*, tasks will be selected and swapped without exceeding the *upper_bound* of each host. Combining these information helps the scheduler make a more balanced task — node mapping than the naïve Spark FIFO scheduling algorithm. We use the OpenCloud trace to evaluate Heuristic MinHeap. Without swapping, the Heuristic MinHeap can achieve a better performance improvement (average 5.7%) than the default Spark FIFO scheduling algorithm (average 2.7%). This is the by-product optimization harvested from shuffle size prediction.

3.2.4 Cope with Multiple Shuffles

Multiple shuffles commonly exist in modern DAG computing. The techniques mentioned in Section 3.2.2 can only handle the ongoing shuffle. For those pending shuffles, it is impossible to predict the sizes. This dilemma can be relieved by having all map tasks of shuffle to be scheduled by DAG framework simultaneously. But doing this introduces large overhead such as extra task serialization. To avoid violating the optimization from framework, we present Accumulate Heuristic Scheduling algorithm (Algorithm 2) to cope with multiple shuffles.

As illustrated in *M_SCHEDULE*, the size of reduce on each node of previously scheduled *shuffles* are counted. When a new shuffle starts, the *M_SCHEDULE* is called to schedule the new one accumulatively. The *size* of each reduce and its corresponding *prob* and *host* in *p_reduces* are updated with data of *shuffles* before *SCHEDULE* is called. When the new task — node mapping is available, if the new *assigned_host* of a reduce does not equal to the original one, the re-shuffle will be triggered to transfer data to a new host for further computing. This re-shuffle is rare since the previously shuffled data in a reduce partition contributes a huge composition. It means in the schedule phase, the *SWAP_TASKS* can help revise the scheduling to match the previous mapping as much as possible while maintaining the good load balance.

4 IMPLEMENTATION

This section presents an overview of the implementation of SCache — an open source cross-framework shuffle data management system with a DAG co-scheduler. In Subsection 4.1, we will first present the system overview, and then we will detail SCache memory management in Subsection 4.2. In Subsection 4.3, we will evaluate cross-framework capability of SCache based on Hadoop Mapreduce and Spark. In the last Subsection 4.4, we will introduce the current status of SCache fault tolerance.

4.1 System Overview

SCache consists of three components: a distributed shuffle data management system, a DAG co-scheduler, and a worker daemon. As a plug-in system, SCache needs to rely on a DAG framework. As shown in Figure 7, SCache employs the legacy master-slaves architecture like GFS [22] for shuffle data management system. The master node of SCache coordinates the shuffle blocks globally with application context. The worker node reserves memory to store blocks. The coordination provides two guarantees: (a) data is stored in memory before tasks start and (b) data is scheduled on-off memory with all-or-nothing and context-aware constraints. The daemon bridges the communication between DAG framework and SCache. The co-scheduler is dedicated to pre-schedule reduce tasks with DAG information and enforce the scheduling results to original scheduler in framework master.

When a DAG job is submitted, the DAG information will be generated in framework master scheduler. Before the computing tasks begin, the shuffle dependencies are determined based on job configuration. The SCache master recognizes all shuffle dependencies in a RPC call as

a shuffle scheduling unit. For each shuffle dependency, the shuffle ID, the type of partitioner, the number of map tasks, and the number of reduce tasks are included. If there is a specialized partitioner, such as range partitioner, in the shuffle dependencies, the daemon will insert a sampling application before the dependent RDDs. We will elaborate the sampling procedure in the Section 4.1.1.

Algorithm 1 Heuristic MinHeap Scheduling for Single Shuffle

```

1: procedure SCHEDULE(m, host_ids, p_reduces)
2:   m ← partition number of map tasks
3:   R ← sort p_reduces by size in descending order
4:   M ← min-heap {host_id → (reduces), size}
5:   idx ← len(R) - 1
6:   while idx ≥ 0 do           ▷ Schedule reduces by MinHeap
7:     M[0].size += R[idx].size
8:     M[0].reduces.append(R[idx])
9:     R[idx].assigned_host ← M[0].host_id
10:    Sift down M[0] by size
11:    idx ← idx - 1
12:  max ← maximum size in M
13:  Convert M to mapping {host_id → (rid_arr), size}
14:  for all reduce in R do           ▷ Heuristic swap by locality
15:    if reduce.assigned_id ≠ reduce.host_id then
16:      p ← reduce.prob
17:      norm ← (p - 1/m) / (1 - 1/m) / 10
18:      upper_bound ← (1 + norm) × max
19:      SWAP_TASKS(M, reduce, upper_bound)
20:  return M
21: procedure SWAP_TASKS(M, reduce, upper_bound)
22:   reduces ← M[reduce.host_id].reduce
23:   candidates ← Select from reduces that assigned_id ≠
24:   host_id and total size closest to reduce.size
25:    $\Delta size \leftarrow sizeOf(candidates) - reduce.size$ 
26:   size_host ← M[reduce.host_id].size -  $\Delta size$ 
27:   size_assigned ← M[reduce.assigned_id].size +  $\Delta size$ 
28:   if size_host ≤ upper_bound and size_assigned ≤ upper_bound then
29:     Swap candidates and reduce
30:     Update size in M
31:     Update assigned_host in candidates and reduce

```

For the hash partitioner, when a map task finishes computing, the SCache daemon process will transfer the map output data from framework worker to the reserved memory through memory copy. After that the slot will be released without blocking on disk operations. When the shuffle map output is received, the SCache worker will notify the master of the block ID and reduce size distribution of this block (see map output in Figure 5). If the collected map output data reach the observation threshold, the DAG co-scheduler will run the scheduling Algorithm 1 and 2 to pre-schedule the reduce tasks and then broadcast the scheduling result to start pre-fetching on each worker. More specifically, when a map task is finished, each node will receive a broadcast message. SCache worker will filter the reduce tasks ID that will be launched on itself and start pre-fetching shuffle data from the remote. In order to force DAG framework to run according to the SCache pre-scheduled results, we insert some lines of codes in framework scheduler. After insertion, DAG scheduler will consult SCache master to get the preferred location for each partition and set *NODE_LOCAL* locality level on corresponding reduce tasks.

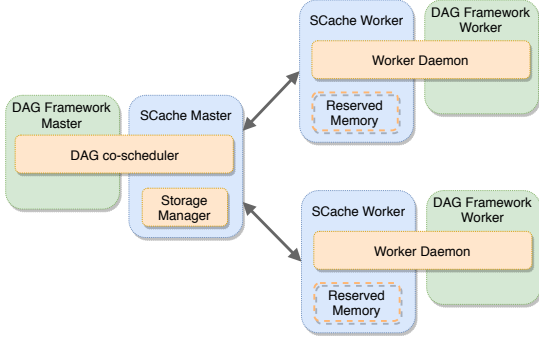


Fig. 7: SCache Architecture

Algorithm 2 Accumulate Heuristic Scheduling for Multi-Shuffles

```

1: procedure M_SCHEDULE( $m, host\_id, p\_reduces, shuffles$ )
2:    $m \leftarrow$  partition number of map tasks  $\triangleright shuffles$  is the
   previous schedule result
3:   for all  $r$  in  $p\_reduces$  do
4:      $r.size \leftarrow shuffles[r.rid].size$ 
5:     if  $shuffles[r.rid].size \geq r.size \times r.prob$  then
6:        $r.prob \leftarrow shuffles[r.rid].size / r.size$ 
7:        $r.host\_id \leftarrow shuffles[r.rid].assigned\_host$ 
8:    $M \leftarrow SCHEDULE(m, host\_id, p\_reduces)$ 
9:   for all  $host\_id$  in  $M$  do  $\triangleright$  Re-shuffle
10:    for all  $r$  in  $M[host\_id].reduces$  do
11:      if  $host \neq shuffles[r.rid].assigned\_host$  then
12:        Re-shuffle data to host
13:       $shuffles[r.rid].assigned\_host \leftarrow host$ 
  return  $M$ 

```

4.1.1 Reservoir Sampling

If the submitted shuffle dependencies contain a range partitioner or a customized non-hash partitioner, the SCache master will send a sampling request to the **framework master**. The daemon then inserts a sampling job before the **submitted computing job**. The sampling job uses a reservoir sampling algorithm [23] on each partition. For the sample number, we set the size to $3 \times \text{number of partitions}$ for balancing overhead and accuracy (it can be tuned by configuration). The sampling job randomly selects some items and performs a local shuffle with partitioner (see Figure 8). At the same time, the items number is counted as the weight. These sampling data will be aggregated by reduce ID on SCache master to predict the reduce partition size. After the prediction, SCache master will call Algorithm 2 and 1 to do the scheduling.

4.2 Memory Management

As mentioned in Section 2.2, though the shuffle size is relatively small, memory management should still be cautious enough to limit the effect of performance of DAG framework. When the size of cached blocks reaches the limitation of reserved memory, SCache flushes some of them to the disk temporarily, and re-fetches them when some cached shuffle blocks are consumed or pre-fetched. To achieve the maximum overall improvement, SCache leverages two constraints to manage the in-memory data — all-or-nothing and context-aware-priority.

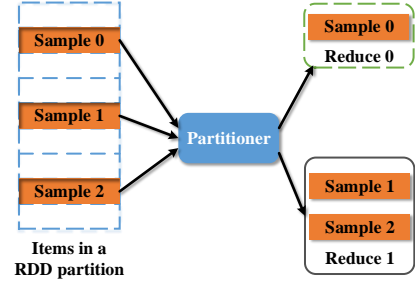


Fig. 8: Reservoir Sampling of One Partition

4.2.1 All-or-Nothing Constraint

Memory cached shuffle blocks can speed up the reduce task execution process. But this acceleration of a single task is necessary but insufficient for a shorter stage completion time. Based on the observation in Section 2.2.4, in most cases one single stage contains multi-rounds of tasks. If one task misses a memory cache and exceeds the original bottleneck of this round, that task might become the new bottleneck and further slow down the whole stage. PACMan [8] has also proved that for multi-round stage/job, the completion time improves in steps when $n \times \text{number of tasks in one round}$ of tasks have data cached simultaneously. Therefore, the cached shuffle blocks need to match the demand of all tasks in one running round at least. We refer to this as the all-or-nothing constraint.

According to all-or-nothing constraint, SCache master leverages the pre-scheduled results to determine the bound of each round, and sets corresponding blocks as the minimum unit of storage management. For those incomplete units, SCache will mark them as the lowest priority.

4.2.2 Context-Aware-Priority Constraint

Unlike the traditional cache replacement schemes such as MIN [24], the cached shuffle data will only be used once (without failure) in DAG computing. That is, the effort of improving hit rate in most legacy schemes can not benefit DAG application, while their approaches can easily violate the all-or-nothing constraint. SCache also leverages application context to select victim storage units when the reserved memory is full.

At first, SCache master searches for the incomplete units and flushes all belonging blocks to disk cluster-widely.

If all the units are completed, SCache selects victims based on two factors — *inter-shuffle units* and *intra-shuffle unit*.

- **Inter-shuffle units:** SCache master follows the scheduling scheme of Spark to determine the inter-shuffle priority. For a FAIR scheduler, Spark balances the resource among task sets, which leads to a higher priority for those having more tasks unfinished. So SCache sets priorities from high to low in a descending order of remaining storage units of a shuffle unit. For a FIFO scheduler, Spark schedules the task set that is submitted first. So SCache sets the priorities according to the submit time of each shuffle unit.
- **Intra-shuffle unit:** SCache also needs to decide the priority among storage units inside a shuffle unit. According to the task scheduling inside a task set of Spark, the

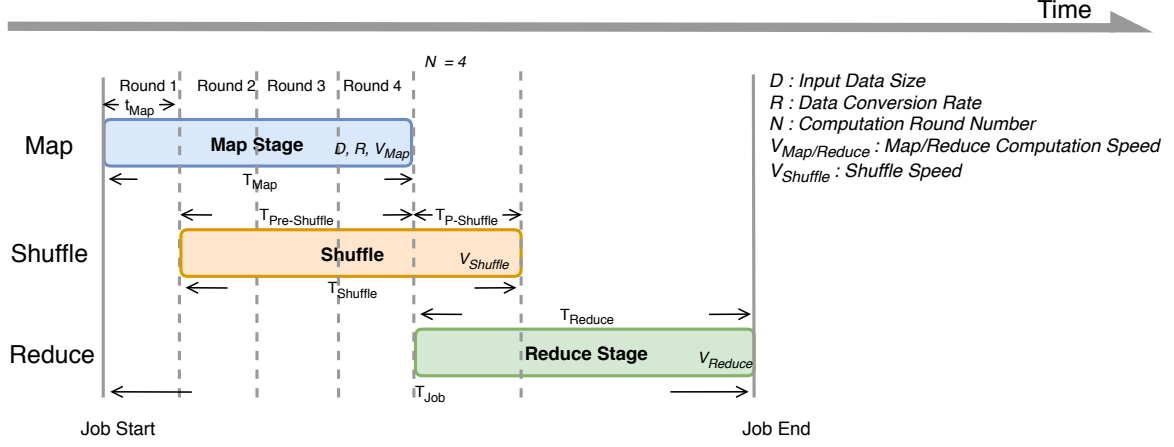


Fig. 9: FRQ Model

tasks with smaller ID will be scheduled firstly. Based on this, SCache can assign the lower priority to storage units with larger task ID.

4.3 Analysis of cross-framework capability

Shuffle optimization of SCache inevitably requires modification on dependent DAG framework. SCache provides API through RPC, such as *putBlock(blockId)*, *getBlock(blockId)*, and *getScheduleResult(shuffleId)*. In order to use SCache, we mainly need to modify the two parts of the framework: (a) We need to insert codes in the DAG scheduler to let framework provide DAG information and follow the order. (b) The shuffle data between the Map task and the Reduce task should be transfer to SCache Storage Management.

To prove the cross-framework capability of SCache, we implement SCache on Hadoop Mapreduce and Spark respectively. Taking Hadoop Mapreduce as an example, we insert codes in ResourceManager and modify codes in MapTask and ReduceTask to call SCache API through RPC. It takes about 400 lines of code. In Spark, We mainly modified the DAGScheduler and the corresponding data fetcher. It only takes about 500 lines of code. Such hundreds of lines of code modification is very small compared to the hundreds of thousands of lines of code in DAG framework. we believe that the costs of enabling SCache on other DAG computing frameworks, such as Tez [3], are also very low.

4.4 Fault tolerance

Since fault tolerance is not a crucial goal of SCache, we have not implemented fault tolerance mechanism. We plan to implement SCache master with Apache ZooKeeper [25] to provide constantly service. For SCache worker, a promising way to prevent failure is to select some backup nodes to store replications of shuffle data during pre-scheduling. In addition, there are also advanced fast recovery techniques such as FineFRC [26]. We leave this to the future work.

5 FRAMEWORK RESOURCES QUANTIFICATION MODEL

In this chapter, we introduce *Framework Resources Quantification*(FRQ) model to describe the performance of DAG

frameworks. FRQ model quantifies computing and I/O resources and displays them in time dimension. According to FRQ model, we can calculate the execution time required by the application under any circumstances, including different DAG framework, hardware environment, and so on. Therefore FRQ model is able to help us analyze the resources scheduling of DAG framework and evaluate their performance. We will first introduce FRQ model in Subsection 5.1. In the following Subsection 5.2, we will use FRQ model to describe three different computation job and analyze their performance. In the last Subsection 5.3, we will use the actual experimental results to verify the FRQ model.

5.1 The FRQ Model

The current distributed parallel computing frameworks mostly use *Directed acyclic graphs*(DAGs) to describe computation logic. A shuffle phase is required between each adjacent DAG computation phase. In order to better analyze the relationship between the computation phase and the shuffle phase, we propose FRQ model. After quantifying computing and I/O resources, FRQ model is able to describe different resource scheduling strategies. For convenience, we introduce FRQ model by taking a simple MapReduce as an example in this section.

Figure 9 shows how the FRQ model describes a MapReduce task. FRQ model has five input parameters:

- 1) Input Data Size(D): The data size of the computation phase.
- 2) Data Conversion Rate(R): The conversion rate of the input data to the shuffle data during a computation phase. This conversion rate depends on the algorithm used in the computation phase.
- 3) Computation Round Number(N): The number of rounds needed to complete the computation phase using the current computation resources. The number of rounds depends on the current computation resources and the settings of the computation job. Take Hadoop MapReduce as an example, suppose we have 50 CPUs and enough memory, the Map phase consists of 200 map tasks. Then we need 4 rounds of computation to complete the Map phase.

- 4) Computation Speed(V_i): The computation speed for each computation phase. The computation speed depends on the algorithm used in the computation phase.
- 5) Shuffle Speed($V_{Shuffle}$): Transmission speed during shuffle. Shuffle speed depends on Network and storage device bandwidth.

We can calculate the execution time of each phase of the job with these five parameters. Obviously, the total execution time of job in this case is the sum of the Map phase time and Reduce phase time:

$$T_{Job} = T_{Map} + T_{Reduce} \quad (3)$$

Map phase time depends on input data size and Map computation speed:

$$T_{Map} = \frac{D}{V_{Map}} \quad (4)$$

The Reduce phase time formula is as follows:

$$T_{Reduce} = \frac{D \times R}{V_{Reduce}} + K \times T_{P_Shuffle} \quad (5)$$

$\frac{D \times R}{V_{Reduce}}$ represents the ideal computation time, and $K \times T_{P_Shuffle}$ represents the calculated overhead. K is empirical value. The overhead depends on the parallel time of shuffle phase and the Reduce phase. The parallel time is denoted by $T_{P_Shuffle}$. And the total time of shuffle phase is represented by $T_{Shuffle}$. Because the computation of the Reduce phase relies on the data transfer results of the Shuffle phase, a portion of the computation in the Reduce phase need to wait for the transfer results. The overhead is caused by these waiting. The FRQ model uses K to indicate the extent of the waiting.

$$T_{Shuffle} = \frac{D}{V_{Shuffle}} \quad (6)$$

For shuffle-heavy computing jobs, we can optimize the job completion time by reducing $T_{P_Shuffle}$. Improving IO speed is a effective way to reduce shuffle time. Another optimization method is to use the idle IO resources in the Map phase for pre-fetching(see Figure 9). Both of the above methods can effectively reduce $T_{P_Shuffle}$. When using FRQ model to describe a computation job, we can easily analyze the resource scheduling strategy of the computation framework. Different computing frameworks may use different resource scheduling strategies. FRQ model can evaluate the scheduling strategies of these computing frameworks and help us optimize them.

5.2 Model Analysis

FRQ model can describe a variety of resource scheduling strategies. First, we analyze a naive scheduling strategy. As shown in Figure 10a, FRQ model describes a Mapreduce job that is fully serially executed. The parallel time of shuffle phase and the Reduce phase is 0, in which case $T_{P_Shuffle}$ is 0. Therefore, the overhead of the Reduce phase is 0. But since shuffle and computation are serial execution, the total execution time of job becomes longer:

$$T_{Job} = T_{Map} + T_{Shuffle} + T_{Reduce} \quad (7)$$

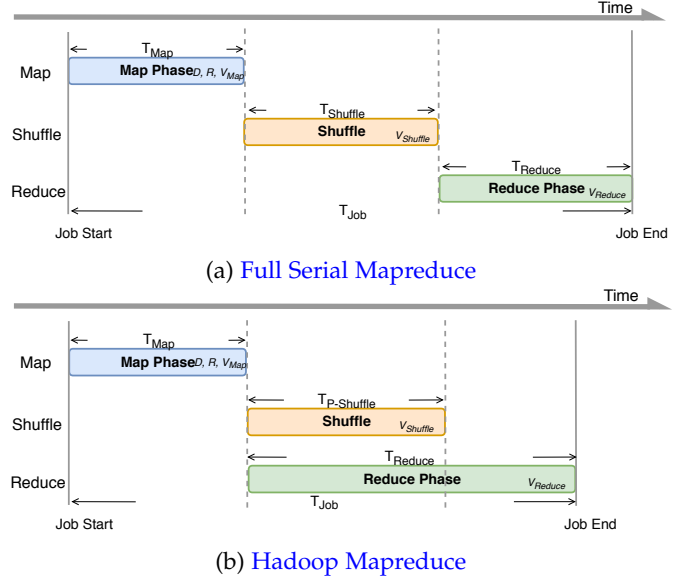


Fig. 10: FRQ Model With Different Scheduling Strategies

Obviously, this is an inefficient scheduling strategy. No computing framework uses this scheduling method. Due to serialization, the IO resource is idle during the Reduce phase and Map phase. The scheduling strategy is naive and has a lot of room for optimization.

Figure 10b shows the scheduling strategy of Hadoop Mapreduce. In Hadoop Mapreduce, Shuffle phase and Reduce phase start at the same time. In this case, $T_{P_Shuffle}$ is equal to $T_{Shuffle}$. Due to the increase in $T_{P_Shuffle}$, the time of Reduce phase will increase(see equation 5). Because the Shuffle phase and the computation phase are executed in parallel, the total execution time of job is the sum of T_{Map} and T_{Reduce} (see equation 3). The execution time of Shuffle phase is hidden in the Reduce phase. This scheduling strategy is much more efficient than the one in Figure 1. However, after analyzing this model, we found that the IO resource in the Map phase is idle. This scheduling strategy can be optimized.

Figure 11 shows the scheduling strategy for Hadoop Mapreduce with SCache (Suppose N is 4). SCache starts pre-fetching and pre-scheduling in the Map phase. This scheduling strategy can make better use of resources and avoid IO resources being idle in the Map phase. According to the design of SCache pre-fetching, we found that using FRQ model to describe the scheduling strategy of SCache needs to distinguish two situations:

- $V_{Map} \times R \geq V_{Shuffle}$ (Figure 11a): The meaning of $V_{Map} \times R$ is the speed at which shuffle data is generated. The meaning of the inequality is that the speed of generating shuffle data($V_{Map} \times R$) is greater than or equal to the shuffle speed($V_{Shuffle}$). When the Round1 of the Map phase ends, the SCache starts shuffling data until the end of the shuffle phase. Due to shuffle speed is slower, the shuffle phase is uninterrupted. SCache transmit the shuffle data generated in the last round of Map phase during the Reduce phase. Therefore $T_{P_Shuffle}$ is equal to one- N of the total time of the

TABLE 1: Hadoop Mapreduce on 4 nodes cluster in FRQ model

	$D: \text{GB}, V_i: \text{GB/s}, T_i: \text{s}$													
	D	R	N	V_{Map}	V_{Reduce}	$V_{Shuffle}$	K	T_{Map}	$T_{Shuffle}$	$T_{P_Shuffle}$	T_{Reduce}	T_{Job}	$ExpT_{Job}$	$Error$
SCache	16	1	2	0.65	1	0.47	0.5	24.62	34.04	21.73	26.87	51.48	55	6.39%
	32	1	4	0.65	1	0.47	0.5	49.23	68.09	31.16	47.58	96.81	104	6.91%
	48	1	6	0.65	1	0.47	0.5	73.85	102.13	40.59	68.29	142.14	151	5.87%
	64	1	8	0.65	1	0.47	0.5	98.46	136.17	50.02	89.01	187.47	193	2.87%
Legacy	16	1	2	0.65	1	0.47	0.6	24.62	34.04	34.04	36.43	61.04	73	16.38%
	32	1	4	0.65	1	0.47	0.6	49.23	68.09	68.09	72.85	122.08	135	9.57%
	48	1	6	0.65	1	0.47	0.6	73.85	102.13	102.13	109.28	183.12	188	2.59%
	64	1	8	0.65	1	0.47	0.6	98.46	136.17	136.17	145.70	244.16	249	1.94%

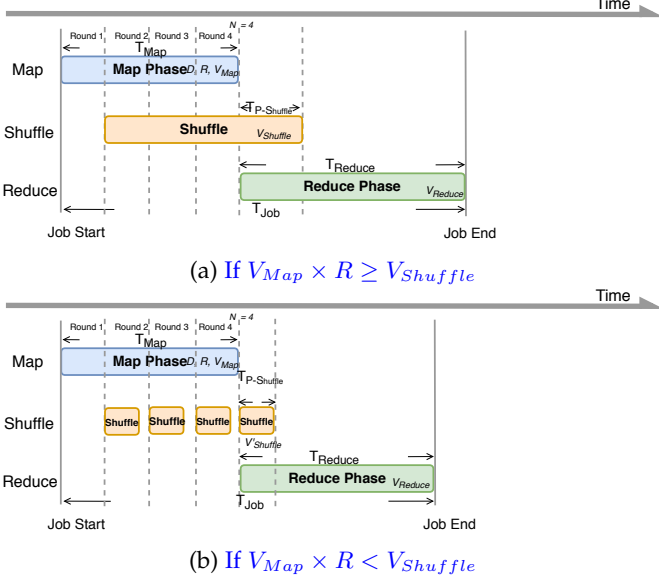


Fig. 11: FRQ Model With SCache

shuffle phase:

$$T_{P_Shuffle} = T_{Shuffle} - \frac{(N-1) \times T_{Map}}{N} \quad (8)$$

- $V_{Map} \times R < V_{Shuffle}$ (Figure 11b): When the speed of generating shuffle data ($V_{Map} \times R$) is less than the shuffle speed ($V_{Shuffle}$), SCache needs to wait for shuffle data to be generated. As Figure 11b shows, the shuffle phase will be interrupted in each Round. Thus $T_{P_Shuffle}$ is equal to the total time of shuffle ($T_{Shuffle}$) minus the time that shuffle is executed in the Map phase:

$$T_{P_Shuffle} = T_{Shuffle} \times \frac{1}{N} \quad (9)$$

Compared to the original Hadoop Mapreduce resource scheduling strategy, Hadoop Mapreduce with SCache reduces $T_{P_Shuffle}$ and thus lessens *Reduce Time* (T_{Reduce}). This is how pre-fetching optimizes the total execution time of job.

5.3 Model Verification

In order to verify FRQ model, we run experiment on two environments. The first environment is on Amazon EC2 and it has 50 m4.xlarge nodes as shown in Section 6.1. Another environment is in our lab. Our lab environment has 4 nodes

and each node has 128GB and 32 CPUs. To simplify the calculation of the FRQ model, we use Hadoop Mapreduce as framework (only three phases: *Map*, *Shuffle*, and *Reduce*) and Terasort as experimental application. We deploy Hadoop with SCache and without SCache on both environments.

Table 1 shows the calculational results of FRQ model in the lab environment. Workload is from 16 GB to 64 GB. D and N are set according to the application parameters. R , V_{Map} , $V_{Shuffle}$, and V_{Reduce} are calculated based on experimental results. K is the empirical value, we set K to 0.5 and 0.6, which reflects that $T_{P_Shuffle}$ has less impact on the Reduce phase in the case of SCache. The formulas of T_{Job} , T_{Map} , T_{Reduce} and $T_{Shuffle}$ are Equation 3, Equation 4, Equation 5 and Equation 6, respectively. In the case of SCache, Terasort on Hadoop Mapreduce satisfies the situation in Figure 11a ($V_{Map} \times R \geq V_{Shuffle}$), thus the formula of $T_{P_Shuffle}$ is Equation 8. In the case of Legacy, since pre-fetching is not used, $T_{P_Shuffle}$ is equal to $T_{Shuffle}$ (see Equation 6). $ExpT_{Job}$ represents the actual experiment data, we calculate *Error* according to T_{Job} and $ExpT_{Job}$. The formula of *Error* is:

$$Error = \frac{ExpT_{Job} - T_{Job}}{T_{Job}} \quad (10)$$

Table 2 shows the calculational results of FRQ model in Amazon EC2 environment. V_{Map} , $V_{Shuffle}$, and V_{Reduce} are modified because of the different hardware devices. We also set K to the same empirical value. The formulas in the table are all the same except $T_{P_Shuffle}$. In this environment, Terasort on Hadoop Mapreduce satisfies the situation in Figure 11b ($V_{Map} \times R < V_{Shuffle}$), thus the formula of $T_{P_Shuffle}$ is Equation 9. In the previous case, $T_{P_Shuffle}$ is still equal to $T_{Shuffle}$.

In order to verify the above-mentioned two cases when using SCache, we monitor network utilization and plot it in Figure 12. Figure 12a shows the utilization of Terasort in the lab environment. The network utilization remains high until shuffle phase is complete. This situation is consistent with Figure 11a. Figure 12b shows the utilization in Amazon EC2 environment. The network utilization is 5 regular peaks, this situation is also consistent with Figure 11b. Therefore, we believe that FRQ model is able to accurately describe framework with SCache.

In terms of accuracy, the experimental values are all greater than the calculated values. This is because the application has some extra overhead at runtime, such as network warm-up, the overhead of allocating slots, and so on. In the case where the input data is small and the total time is short,

TABLE 2: Hadoop Mapreduce on 50 AWS m4.xlarge nodes cluster in FRQ model

	$D: \text{GB}, V_i: \text{GB/s}, T_i: \text{s}$													
	D	R	N	V_{Map}	V_{Reduce}	$V_{Shuffle}$	K	T_{Map}	$T_{Shuffle}$	$T_{P_Shuffle}$	T_{Reduce}	T_{Job}	$ExpT_{Job}$	$Error$
SCache	128	1	5	1.15	1.46	1.4	0.5	111.30	91.43	18.29	96.81	208.12	232	10.29%
	256	1	5	1.15	1.46	1.4	0.5	222.61	182.86	36.57	193.63	416.24	432	3.65%
	384	1	5	1.15	1.46	1.4	0.5	333.91	274.29	54.86	290.44	624.36	685	8.85%
Legacy	128	1	5	1.15	1.46	1.4	0.6	111.30	91.43	91.43	142.53	253.83	266	4.57%
	256	1	5	1.15	1.46	1.4	0.6	222.61	182.86	182.86	285.06	507.67	524	3.12%
	384	1	5	1.15	1.46	1.4	0.6	333.91	274.29	274.29	427.59	761.50	776	1.87%

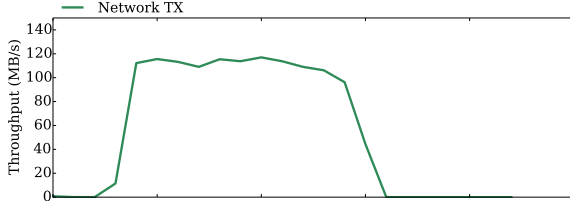
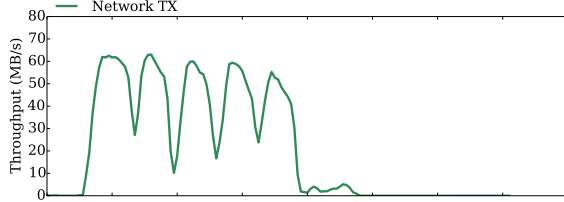
(a) If $V_{Map} \times R \geq V_{Shuffle}$ (b) If $V_{Map} \times R < V_{Shuffle}$

Fig. 12: Network utilization on Hadoop Mapreduce with SCache

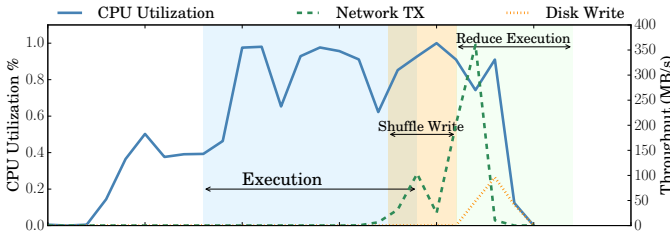


Fig. 13: CPU utilization and I/O throughput of a node during a Spark single shuffle application With SCache

the error caused by the overhead is amplified. Overall, the error between T_{Job} and $ExpT_{Job}$ is basically below 10%, such errors are within tolerance. Therefore, We believe that FRQ model can accurately describe DAG framework.

6 EVALUATION

This section reveals the evaluation of SCache with comprehensive workloads and benchmarks. First we run a simple DAG job with single shuffle to analyze hardware utilization and impact of shuffle optimization from the scope of a task to a job. Then we use a recognized shuffle intensive benchmark — Terasort to evaluate SCache with different data partition schemes.

In order to prove the performance gain of SCache with a real production workload, we also evaluate Spark TPC-

DS¹ and present the overall performance improvement. To prove SCache compatibility as a cross-framework plugin, we implemented SCache on both Hadoop MapReduce and Spark. Due to the simple DAG computing in Hadoop MapReduce, we only use Terasort as a shuffle-heavy benchmark to evaluate the performance of Hadoop Mapreduce with SCache. At last, we measure the overhead of weighted reservoir sampling. In short, SCache can decrease 89% time of Spark shuffle. Furthermore, SCache achieves a 75% and 50% improvement of reduce stage completion time respectively in simple DAG application and Terasort without introducing extra network traffic. More impressively, the overall completion time of TPC-DS can be improved 40% on average by applying optimization from SCache. Meanwhile, Hadoop Mapreduce with Scache optimize job completion time by up to 15% and an average of 13%

6.1 Setup

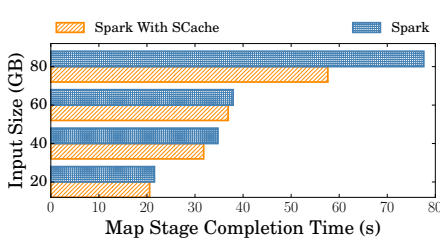
We implement a Spark demon to enable shuffle optimization as a representative. We run our experiments on a 50 m4.xlarge nodes cluster on Amazon EC2². Each node has 16GB memory and 4 CPUs. The network bandwidth is not specifically provided by Amazon. Our evaluations reveal the bandwidth is about 300 Mbps (see Figure 2).

6.2 Simple DAG Analysis

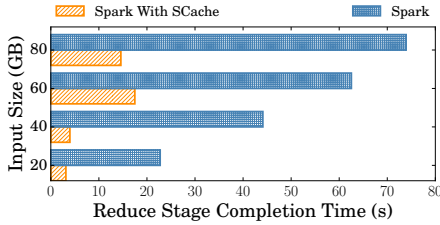
We first run the same single shuffle test as mentioned in Figure 2. As shown in Figure 13, the hardware utilization is captured from one node during the job. Note that since the completion time of whole job is about 50% less than Spark without SCache, the duration of Figure 13 is cut in half as well. An overlap among CPU, disk, and network can be easily observed in Figure 13. That is, SCache prevents the computing process from being cutting off by the I/O operations with a fine-grained resource allocation. It ensures the overall CPU utilization of the cluster stays in a high level. In addition, the decoupling of shuffle write helps free the CPU earlier and leads to a faster map task computation. At the same time, the shuffle pre-fetching in the early stage of map phase shifts the network transfer completion time, which decreases reduce stage completion time. The combination is the main performance gain we achieved on the scope of hardware utilization by SCache.

The performance evaluation in Figure 15 shows the consistent results with our analysis on hardware utilization. For each stage, we run 5 rounds of tasks with different input

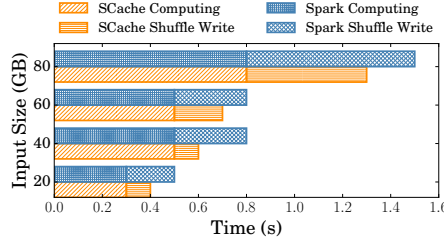
1. <https://github.com/databricks/spark-sql-perf>
2. <http://aws.amazon.com/ec2/>



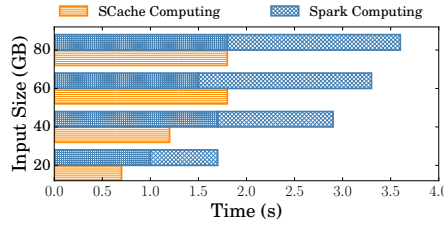
(a) Map Stage Completion Time



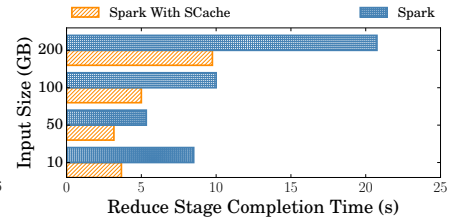
(b) Reduce Stage Completion Time



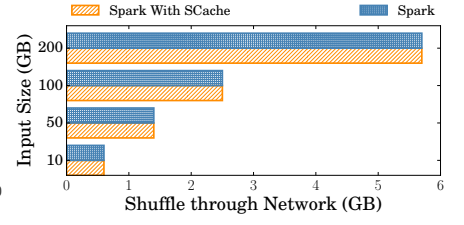
(a) Median Task in Map Stages



(b) Median Task in Reduce Stages



(a) Reduce Stage of First Shuffle



(b) Network Traffic of Second Shuffle

Fig. 14: Stage Completion Time of Single Shuffle Test

Fig. 15: Median Task Completion Time of Single Shuffle Test

Fig. 16: Terasort Evaluation

size. By running spark with SCache, the completion time of map stage can be reduced 10% on average. For reduce stage, instead, SCache achieves a 75% performance gain in the completion time of the reduce stage.

A detail analysis into the nutshell of varied overall performance gain on different stages is presented with Figure 15. For each stage, we pick the median task. About 40% of shuffle write time can be eliminated by SCache (Figure 15a) in a map task. Because the serialization of data is CPU intensive [14] and it is inevitable while moving data out of JVM, SCache can not eliminate the whole phase of shuffle write. It results in a less performance gain in the map stage. On the other hand, the network transfer introduces a significantly latency in shuffle read of reduce task (Figure 15b). By doing shuffle data pre-fetching for the reduce tasks, the shuffle read time decreases 100%. That is, the explicit network transfer is perfectly overlapped by SCache shuffle pre-fetching. In overall, SCache can help Spark decrease by 89% time in the whole shuffle process. In addition, heuristic reduce tasks scheduling achieves better load balance in cluster than the Spark default FIFO scheduling which may randomly assign two heavy tasks on a single node. So that we can have a significant performance gain in the completion time of the reduce stage.

6.3 Terasort

We also evaluate Terasort — a recognized shuffle intensive benchmark for distributed system analysis. Terasort consists of two consecutive shuffles. The first shuffle reads the input data and uses a customized hash partition function for repartitioning. The second shuffle partitions the data through a range partitioner. As the range bounds set by range partitioner almost match the same pattern of the first shuffle, almost 93% of input data is from one particular map task for each reduce task. So we take the second shuffle as an extreme case to evaluate the scheduling locality for SCache.

As shown in Figure 16a, we present the first shuffle as the evaluation of shuffle optimization. At the same time,

we use the second shuffle to evaluate in the dimension of scheduling locality (Figure 16b). For the first shuffle, Spark with SCache runs $2 \times$ faster during the reduce stage with the input data in a range from 10GB to 200GB. At the same time, Figure 16b reveals that SCache pre-scheduling produces exactly same network traffic of the second shuffle as Spark, which implies that SCache pre-scheduling can obtain the best locality while balancing the load. In contrast, Spark delays scheduling reduce tasks with the shuffle map output to achieve this optimum.

6.4 Hadoop Mapreduce with SCache

To prove SCache compatibility as a cross-framework plugin, we also implemented SCache on Hadoop MapReduce as the external shuffle service and co-scheduler. Although *pre-scheduling reduce tasks* is not critical for the simple DAG computing in Hadoop MapReduce, some shuffle-heavy jobs on Hadoop Mapreduce can still be optimized by SCache shuffle data management.

Figure 19 shows the hardware resource utilization of Hadoop Mapreduce running Terasort. Both figures have the same proportion of time. Hadoop MapReduce with SCache brings 15% of total time optimization with 384GB input data size. As shown in the Figure 19b, Hadoop MapReduce without SCache writes intermediate data locally in the Map phase. The shuffle phase and reduce phase start simultaneously. Because the large amount of shuffle data reaches the network bottleneck, The beginning part of reduce needs to wait needs to wait for network transfer. This causes the CPU resources to be idle. As shown in the Figure 19a, Hadoop Mapreduce with SCache start pre-fetching in the Map phase. This avoids the Reduce phase waiting for the shuffle data. Furthermore, pre-fetching utilize the idle IO throughput in the Map phase. As shown in Figure 20, after better fine-grained utilization of hardware resources, Hadoop Mapreduce with SCache optimize Terasort overall completion time by up to 15% and an average of 13% with input data sizes from 128GB to 512GB.

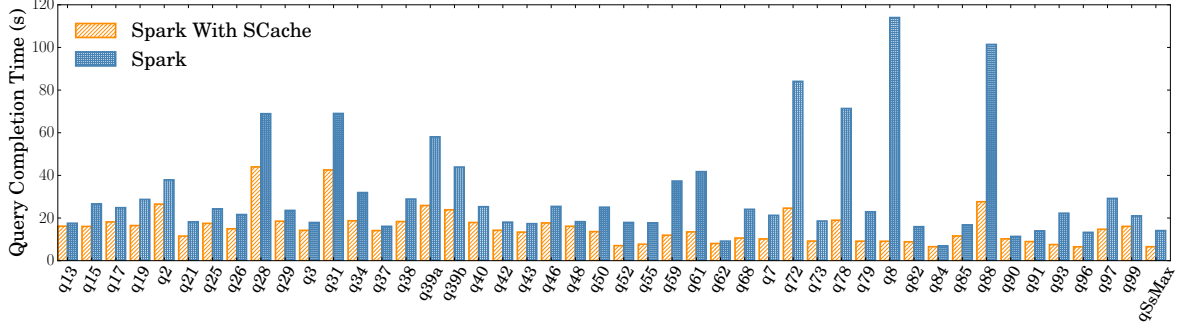


Fig. 17: TPC-DS Benchmark Evaluation

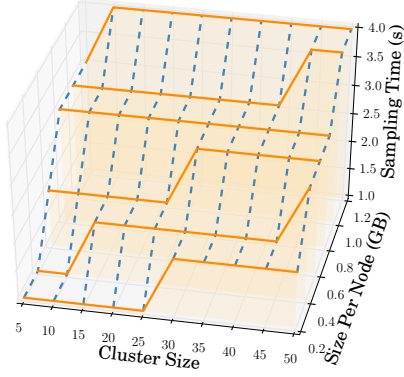
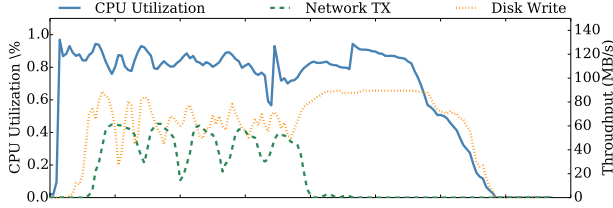
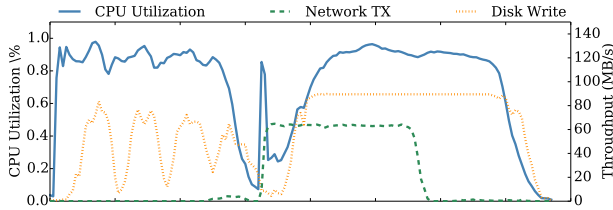


Fig. 18: Sampling Overhead



(a) Hadoop Mapreduce With SCache



(b) Hadoop Mapreduce Without SCache

Fig. 19: CPU utilization and I/O throughput of a node during a Hadoop Mapreduce Terasort job

6.5 Production Workload

We also evaluate some shuffle heavy queries from TPC-DS. TPC-DS benchmark is designed for modeling multiple users submitting varied queries (e.g. ad-hoc, interactive OLAP, data mining, etc.). TPC-DS contains 99 queries and is considered as the standardized industry benchmark for testing big data systems. We evaluate the performance of Spark with SCache by picking some of the TPC-DS queries with shuffle intensive attribute. As shown in Figure 17, the horizontal axis is query number, and the vertical axis

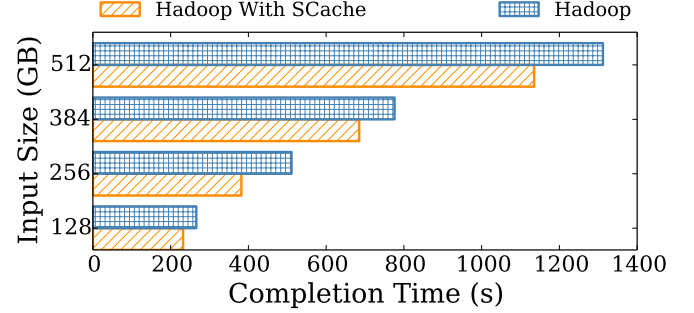


Fig. 20: Hadoop Mapreduce Terasort completion time

is query completion time. Spark with SCache outperforms the original Spark in almost all the queries. Furthermore, in many queries, Spark with SCache outperforms original Spark by an order of magnitude. The overall reduction portion of query time that SCache achieved is 40% on average. Since this evaluation presents the overall job completion time of queries, we believe that our shuffle optimization is promising.

6.6 Overhead of Sampling

In this part, we evaluate the overhead of sampling with different input data sizes on one node and cluster scales. In Figure 18, the overhead of sampling only grows with the increase of input size on each node. But it remains relatively stable when the cluster size scales up. It makes SCache a scalable system in cluster.

7 RELATED WORK

This section describes related work of shuffle optimization.

Modeling: Most well-known DAG computing frameworks use similar *Bulk Synchronize Parallel*(BSP) [27] model to control data synchronization in each computing phase, i.e. *stage* in Spark, *superstep* in Pregel [28] and so on (Apache Hadoop Mapreduce¹ can also be considered as a special case of only one superstep in BSP model). In the process of optimizing the shuffle phases between adjacent computing phases, we design a performance model to assist in analyzing computing process. Inspired by [29], [30], and [31], we present *Framework Resources Quantification*(FRQ) model to describe the performance of DAG frameworks.

1. <http://hadoop.apache.com/>

In [29], the author designs a model and divides execution into three parts: Map, Shuffle and Reduce. And then the author uses a greedy algorithm to roughly calculate the max, min and mean execution time of Map and Reduce phases. In [30], the author presents a detailed set of mathematical performance models for describing the execution of a MapReduce job. The execution is separated into the phases: Read, Map, Collect, Spill, Merge, Shuffle, Merge, and Reduce. The performance models describe each above-mentioned phases and combine into a overall Mapreduce job model. However, none of the above models meet our requirements for analyzing shuffle. The model in [29] is not able to accurately describe the overhead caused by shuffle under different scheduling strategies. The model in [30] calculates overhead of various phases including shuffle, but most of them are redundant. Different from these models, our FRQ model quantifies computing and I/O resources and displays them in time dimension. FRQ model focuses on describing the overhead caused by shuffle in different scheduling strategies, which satisfy our demand.

Pre-scheduling: Slow-start from Apache Hadoop MapReduce is the most classic approach to handle shuffle overhead. Starfish [32] gets sampled data statistics for self-tuning system parameters (e.g. slow-start, map and reduce slot ratio, etc). DynMR [33] dynamically starts reduce tasks in late map stage to decrease the time of waiting for outputs of map tasks. All of them left the explicit I/O time in an occupied computation slot. Instead, SCache can start shuffle pre-fetching without consuming slots. iShuffle [20] decouples shuffle from reducers and designs a centralized shuffle controller. But it can neither handle multiple shuffles nor schedule multiple rounds of reduce tasks. iHadoop [34] aggressively pre-schedules tasks in multiple successive stages, in order to start fetching data from previous stage earlier. But we have proved that randomly assign tasks may hurt the overall performance in Section 3.2.1. Different from these works, SCache pre-schedules multiple shuffles without breaking load balancing by combining DAG information and heuristic algorithms.

Delay-scheduling: Delay Scheduling [10] delays tasks assignment to get better data locality, which can reduce the network traffic. ShuffleWatcher [13] delays shuffle fetching when network is saturated. At the same time, it achieves better data locality. Both Quincy [9] and Fair Scheduling [35] can reduce shuffle data by optimizing data locality of map tasks. But all of them can not mitigate explicit I/O in both map and reduce tasks. In addition their optimizations fluctuate under different network performances and data distributions, whereas SCache can provide a stable performance gain by shuffle in-memory caching and pre-fetching.

Network layer optimization: Varys [36] and Aalo [37] provide the network layer optimization for shuffle transfer. Though the efforts are limited throughout whole shuffle process, they can be easily applied on SCache to further improve the performance.

8 CONCLUSION

In this paper, we present SCache, a cross-framework open source shuffle optimization system for DAG computing frameworks. SCache decouples the shuffle from computing

pipeline and leverages shuffle data pre-fetching to mitigate the I/O overhead of the whole system. By scheduling tasks with application context, SCache bridges the gap among computing stages. At the expense of inserting hundred lines of code, we adapt Spark and Hadoop Mapreduce to SCache. Our evaluation result shows that SCache can provide a promising speedup to the DAG framework. Furthermore, we propose *Framework Resources Quantification*(FRQ) model to assist in analyzing shuffle process of DAG computing frameworks. At last we use FRQ model to verify SCache shuffle optimization by mathematics. Therefore we believe that SCache is a simple and efficient system to enhance the performance of most DAG computing frameworks.

9 ACKNOWLEDGE

This work was supported in part by National Key Research & Development Program of China (No. 2016YFB1000502), National NSF of China (NO. 61672344, 61525204, 61732010), and Shanghai Key Laboratory of Scalable Computing and Systems.

REFERENCES

- [1] M. Isard, M. Budiu, Y. Yu, A. Birrell, and D. Fetterly, "Dryad: distributed data-parallel programs from sequential building blocks," in *ACM SIGOPS operating systems review*, vol. 41, no. 3. ACM, 2007, pp. 59–72.
- [2] 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.
- [3] B. Saha, H. Shah, S. Seth, G. Vijayaraghavan, A. Murthy, and C. Curino, "Apache tez: A unifying framework for modeling and building data processing applications," in *Proceedings of the 2015 ACM SIGMOD international conference on Management of Data*. ACM, 2015, pp. 1357–1369.
- [4] M. Chowdhury, M. Zaharia, J. Ma, M. I. Jordan, and I. Stoica, "Managing data transfers in computer clusters with orchestra," in *ACM SIGCOMM Computer Communication Review*, vol. 41, no. 4. ACM, 2011, pp. 98–109.
- [5] C. Xie, R. Chen, H. Guan, B. Zang, and H. Chen, "Sync or async: Time to fuse for distributed graph-parallel computation," in *Proceedings of the 20th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPOPP 2015. New York, NY, USA: ACM, 2015, pp. 194–204. [Online]. Available: <http://doi.acm.org/10.1145/2688500.2688508>
- [6] S. Babu, "Towards automatic optimization of mapreduce programs," in *Proceedings of the 1st ACM Symposium on Cloud Computing*, ser. SoCC '10. New York, NY, USA: ACM, 2010, pp. 137–142. [Online]. Available: <http://doi.acm.org/10.1145/1807128.1807150>
- [7] H. Li, A. Ghodsi, M. Zaharia, S. Shenker, and I. Stoica, "Tachyon: Reliable, memory speed storage for cluster computing frameworks," in *Proceedings of the ACM Symposium on Cloud Computing*. ACM, 2014, pp. 1–15.
- [8] G. Ananthanarayanan, A. Ghodsi, A. Wang, D. Borthakur, S. Kandula, S. Shenker, and I. Stoica, "Pacman: Coordinated memory caching for parallel jobs," in *Proceedings of the 9th USENIX conference on Networked Systems Design and Implementation*. USENIX Association, 2012, pp. 20–20.
- [9] M. Isard, V. Prabhakaran, J. Currey, U. Wieder, K. Talwar, and A. Goldberg, "Quincy: fair scheduling for distributed computing clusters," in *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles*. ACM, 2009, pp. 261–276.
- [10] M. Zaharia, D. Borthakur, J. Sen Sarma, K. Elmeleegy, S. Shenker, and I. Stoica, "Delay scheduling: a simple technique for achieving locality and fairness in cluster scheduling," in *Proceedings of the 5th European conference on Computer systems*. ACM, 2010, pp. 265–278.

- [11] J. Dean and S. Ghemawat, "Mapreduce: simplified data processing on large clusters," *Communications of the ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- [12] M. Zaharia, R. S. Xin, P. Wendell, T. Das, M. Armbrust, A. Dave, X. Meng, J. Rosen, S. Venkataraman, M. J. Franklin, A. Ghodsi, J. Gonzalez, S. Shenker, and I. Stoica, "Apache spark: A unified engine for big data processing," *Commun. ACM*, vol. 59, no. 11, pp. 56–65, Oct. 2016. [Online]. Available: <http://doi.acm.org/10.1145/2934664>
- [13] F. Ahmad, S. T. Chakradhar, A. Raghunathan, and T. Vijaykumar, "Shufflewatcher: Shuffle-aware scheduling in multi-tenant mapreduce clusters," in *USENIX Annual Technical Conference*, 2014, pp. 1–12.
- [14] K. Ousterhout, R. Rasti, S. Ratnasamy, S. Shenker, B.-G. Chun, and V. ICSL, "Making sense of performance in data analytics frameworks," in *NSDI*, vol. 15, 2015, pp. 293–307.
- [15] B. Fitzpatrick, "Distributed caching with memcached," *Linux J.*, vol. 2004, no. 124, pp. 5–, Aug. 2004. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1012889.1012894>
- [16] D. Ongaro, S. M. Rumble, R. Stutsman, J. Ousterhout, and M. Rosenblum, "Fast crash recovery in ramcloud," in *Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles*. ACM, 2011, pp. 29–41.
- [17] Y. Kwon, M. Balazinska, B. Howe, and J. Rolia, "Skewtune: mitigating skew in mapreduce applications," in *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*. ACM, 2012, pp. 25–36.
- [18] G. Ananthanarayanan, S. Kandula, A. G. Greenberg, I. Stoica, Y. Lu, B. Saha, and E. Harris, "Reining in the outliers in mapreduce clusters using mantri," in *OSDI*, vol. 10, no. 1, 2010, p. 24.
- [19] B. Gufler, N. Augsten, A. Reiser, and A. Kemper, "Load balancing in mapreduce based on scalable cardinality estimates," in *Data Engineering (ICDE), 2012 IEEE 28th International Conference on*. IEEE, 2012, pp. 522–533.
- [20] D. Cheng, J. Rao, Y. Guo, and X. Zhou, "Improving mapreduce performance in heterogeneous environments with adaptive task tuning," in *Proceedings of the 15th International Middleware Conference*. ACM, 2014, pp. 97–108.
- [21] A. Verma, L. Cherkasova, and R. H. Campbell, "Resource provisioning framework for mapreduce jobs with performance goals," in *ACM/IFIP/USENIX International Conference on Distributed Systems Platforms and Open Distributed Processing*. Springer, 2011, pp. 165–186.
- [22] S. Ghemawat, H. Gobioff, and S.-T. Leung, "The google file system," in *ACM SIGOPS operating systems review*, vol. 37, no. 5. ACM, 2003, pp. 29–43.
- [23] J. S. Vitter, "Random sampling with a reservoir," *ACM Transactions on Mathematical Software (TOMS)*, vol. 11, no. 1, pp. 37–57, 1985.
- [24] L. A. Belady, "A study of replacement algorithms for a virtual-storage computer," *IBM Systems journal*, vol. 5, no. 2, pp. 78–101, 1966.
- [25] P. Hunt, M. Konar, F. P. Junqueira, and B. Reed, "Zookeeper: Wait-free coordination for internet-scale systems," in *USENIX annual technical conference*, vol. 8, 2010, p. 9.
- [26] R. Christodouloupoulou, K. Manassiev, A. Bilas, and C. Amza, "Fast and transparent recovery for continuous availability of cluster-based servers," in *Proceedings of the Eleventh ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPOPP '06. New York, NY, USA: ACM, 2006, pp. 221–229. [Online]. Available: <http://doi.acm.org/10.1145/1122971.1123005>
- [27] L. G. Valiant, "A bridging model for parallel computation," *Communications of the ACM*, vol. 33, no. 8, pp. 103–111, 1990.
- [28] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski, "Pregel: a system for large-scale graph processing," in *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data*. ACM, 2010, pp. 135–146.
- [29] A. Verma, L. Cherkasova, and R. H. Campbell, "Aria: automatic resource inference and allocation for mapreduce environments," in *Proceedings of the 8th ACM international conference on Autonomic computing*. ACM, 2011, pp. 235–244.
- [30] H. Herodotou, "Hadoop performance models," *arXiv preprint arXiv:1106.0940*, 2011.
- [31] J. Polo, D. Carrera, Y. Becerra, J. Torres, E. Ayguadé, M. Steinder, and I. Whalley, "Performance-driven task co-scheduling for mapreduce environments," in *2010 IEEE Network Operations and Management Symposium-NOMS 2010*. IEEE, 2010, pp. 373–380.
- [32] H. Herodotou, H. Lim, G. Luo, N. Borisov, L. Dong, F. B. Cetin, and S. Babu, "Starfish: A self-tuning system for big data analytics," in *Cidr*, vol. 11, no. 2011, 2011, pp. 261–272.
- [33] 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.
- [34] E. Elnikety, T. Elsayed, and H. E. Ramadan, "ihadoop: asynchronous iterations for mapreduce," in *Cloud Computing Technology and Science (CloudCom), 2011 IEEE Third International Conference on*. IEEE, 2011, pp. 81–90.
- [35] Y. Wang, J. Tan, W. Yu, L. Zhang, X. Meng, and X. Li, "Preemptive reducetask scheduling for fair and fast job completion," in *ICAC*, 2013, pp. 279–289.
- [36] M. Chowdhury, Y. Zhong, and I. Stoica, "Efficient coflow scheduling with varys," in *ACM SIGCOMM Computer Communication Review*, vol. 44, no. 4. ACM, 2014, pp. 443–454.
- [37] M. Chowdhury and I. Stoica, "Efficient coflow scheduling without prior knowledge," in *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication*, ser. SIGCOMM '15. New York, NY, USA: ACM, 2015, pp. 393–406. [Online]. Available: <http://doi.acm.org/10.1145/2785956.2787480>

Rui Ren

PLACE
PHOTO
HERE

Chunghsuan Wu

PLACE
PHOTO
HERE

Zhouwang Fu

PLACE
PHOTO
HERE



Tao Song Tao Song is currently working as a postdoc at Shanghai Jiao Tong University in China. He received his Ph.D. degree in computer science and M.Eng. degree in software engineering from Shanghai Jiao Tong University. His research interests include data center networking, cloud computing, artificial intelligence and swarm intelligence.



Zhengwei Qi Zhengwei Qi received the BEng and MEng degrees from Northwestern Polytechnical University, in 1999 and 2002, and the PhD degree from Shanghai Jiao Tong University, in 2005. Currently, he is a professor in the School of Software, Shanghai Jiao Tong University (China). His research interests include distributed computing, virtualized security, model checking, program analysis and embedded systems.



Haibing Guan Haibing Guan received the PhD degree from Tongji University, in 1999. He is a professor of School of Electronic, Information and Electronic Engineering, Shanghai Jiao Tong University, and the director of the Shanghai Key Laboratory of Scalable Computing and Systems. His research interests include distributed computing, network security, network storage, green IT, and cloud computing.