

The Misbelief in Delay Scheduling

Derek Schatzlein
Purdue University

Srivatsan Ravi
Purdue University

Youngtae Noh
Inha University

Masoud
Saeida Ardekani^{*}
Purdue University

Patrick Eugster[†]
Purdue University, TU
Darmstadt

ABSTRACT

Big-data processing frameworks like Hadoop and Spark, often used in multi-user environments, have struggled to achieve a balance between the full utilization of cluster resources and fairness between users. In particular, data locality becomes a concern, as enforcing fairness policies may cause poor placement of tasks in relation to the data on which they operate. To combat this, the schedulers in many frameworks use a heuristic called delay scheduling, which involves waiting for a short, constant interval for data-local task slots to become free if none are available; however, a fixed delay interval is inefficient, as the ideal time to delay varies depending on input data size, network conditions, and other factors.

We propose an adaptive solution (Dynamic Delay Scheduling), which uses a simple feedback metric from finished tasks to adapt the delay scheduling interval for subsequent tasks at runtime. We present a dynamic delay implementation in Spark, and show that it outperforms a fixed delay in TPC-H benchmarks. Our preliminary experiments confirm our intuition that job latency in batch-processing scheduling can be improved using simple adaptive techniques with almost no extra state overhead.

1. INTRODUCTION

Big-data processing frameworks like Hadoop [5], and Spark [9] have been widely used in many companies such as Google, Yahoo, Facebook, and Microsoft. Their popularity brings forth new problems in the realm of task scheduling and management in these systems, particularly with the rise of the cloud as a platform to host these frameworks.

In this paper, we explore the problem of sharing a cluster between users while preserving the efficiency of systems

^{*}Now at Samsung Research America.

[†]Funded by the German Research Foundation under project MAKI, and by European Research Council under project Lightweight Verification of Software

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

DCC'16, July 25-28, 2016, Chicago, IL, USA

© 2016 ACM. ISBN 978-1-4503-4220-9/16/07...\$15.00

DOI: <http://dx.doi.org/10.1145/2955193.2955203>

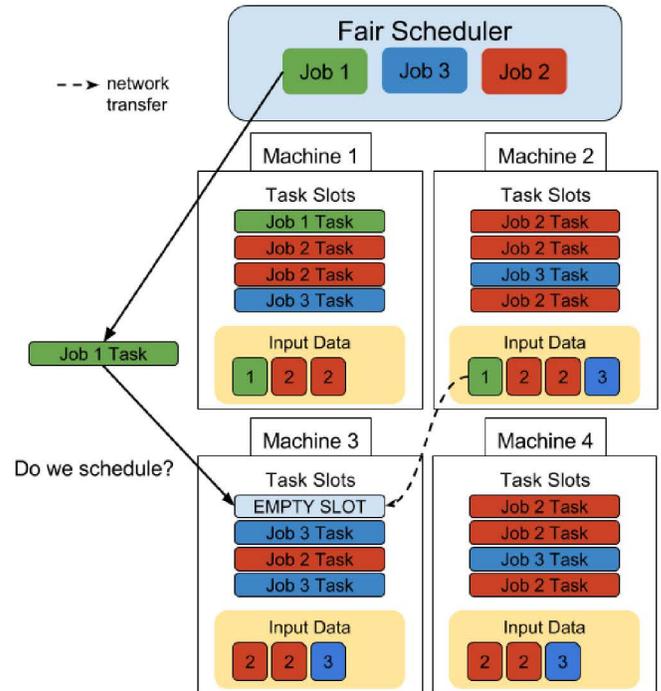


Figure 1: A compute cluster, with a fair scheduler that must fill an empty slot, considering both fairness and data locality.

like Hadoop and Spark — specifically, finding a trade-off between locality (i.e., the placement of computation near its input data) and fairness. Locality is important for task (and thus job) performance, as it reduces the amount of data that needs to travel across the network, while fairness is important for clusters to be useful to multiple users at load.

Figure 1 demonstrates these tradeoffs between data locality and fairness, in a simple four-machine cluster. We consider a scenario where a “job” consists of many smaller tasks, and each node in the cluster has a finite number of identical “slots” which provide resources for tasks. Every task has input data, located in some distributed file system in blocks across the cluster. We assume that the latency of transferring these blocks across the network is high. When a slot becomes available (a slot representing the resources to run a single task), a scheduler that provides fair sharing needs to schedule tasks from the job that is farthest below its fair share first. In Figure 1, the job furthest below the

fair share is clearly Job 1, so a task from Job 1 will be chosen for the empty slot; however, Job 1 has no input data on Machine 3, meaning that the task will need to transfer a block of data over the network in order to proceed, incurring additional overhead and causing worse performance. The question becomes: what should a fair scheduler do in this situation? Should it sacrifice performance to enforce fairness, or choose another task with local data and remedy the fairness imbalance later?

The Hadoop Fair Scheduler (HFS) [7] is an influential example of a cluster framework scheduler that seeks to find this balance. The HFS was designed with two main goals: *Fair sharing*: dividing resources using max-min fair sharing [11] to achieve statistical multiplexing, and *Data locality*: placing computations (tasks) near their input data as often as possible, to maximize system throughput.

To achieve the first goal (fair sharing), a scheduler must reallocate resources between jobs when the number of jobs changes; however, a strict implementation of fair sharing compromises locality, because the job to be scheduled next according to fairness might not have data on the nodes that are currently free. To resolve this problem, the HFS relaxes fairness slightly through a simple algorithm called *delay scheduling* [8], in which a job waits for a fixed amount of time for additional scheduling opportunities on nodes that have data for it, if none are available.

The original delay scheduling work [8] reported that a small amount of waiting is enough to bring locality close to 100% (the chosen default value is three seconds). Since its inception, delay scheduling has been used in other cluster framework schedulers, such as Spark, as well as in resource managers, such as YARN [6] and Mesos [2].

Delay scheduling provides a simple solution for the scenario in Figure 1. Job 1 will simply be delayed for a fixed interval (and a task from Job 3 scheduled in its place) in the hopes that data-local slots will open up in the future; however, delay scheduling assumes that a fixed waiting time always results in a good probability of achieving data locality, and that it will always benefit task completion time. In reality, the decision is not so simple, as the frequency of slot openings, and the penalty incurred from missing data locality, changes based on a number of factors, including network conditions, input data size, disk i/o usage, and others. These conditions may change over time, or even during the execution of a job. As well, with clusters in the public cloud becoming so common, cluster conditions can change as simultaneous user load changes, or after instances are stopped and restarted, potentially in a different physical configuration. So, the decision of whether or not to delay (and for how long) can not be efficiently solved with a fixed, unchanging delay interval. We propose that an *adaptive* mechanism is needed, in order to properly judge and exploit these trade-offs between fairness and locality, based on feedback from the system.

In this paper, we present Dynamic Delay Scheduling, a *simple, low-state*, adaptive solution for task scheduling in the cloud, which improves job latency over fixed-interval delay scheduling. Dynamic Delay Scheduling operates by using feedback from finished tasks in the system, which report their perceived overhead from remote data reads. Our scheduler then adapts the delay scheduling interval to be used to schedule subsequent tasks, using a running average of the overhead metrics reported as the interval itself. In

this way, our adaptive solution minimizes completion time by finding the shifting balance between: 1) immediate fairness policy enforcement, and 2) waiting for data locality, based on the severity of network overhead, for increased performance. This paper contributes an outline of this adaptive solution. We have also implemented a simple prototype in Spark. Preliminary experiments demonstrate that an adaptive delay outperforms the current Spark distribution (which uses a fixed delay scheduling interval) in TPC-H workloads.

The remainder of this paper is structured in the following way. Section 2 presents background on fixed delay scheduling as a concept. Section 3 discusses how fixed delay scheduling is implemented/mapped to real systems. Section 4 defines the terminology and design of Dynamic Delay Scheduling, our simple adaptive scheduling solution. Section 5 describes our lightweight implementation in Apache Spark. Section 6 demonstrates our evaluation of Dynamic Delay Scheduling, using our implementation in Spark under several workloads. Finally, Section 7 overviews related work while Section 8 provides closing remarks, and well as some future direction for further applying adaptive algorithms to other facets of task scheduling.

2. DELAY SCHEDULING BACKGROUND

Delay scheduling is a heuristic designed to improve the data locality of tasks in the context of fairness policy enforcement. We consider a cluster computing setting, where each application (or job), when submitted, is divided into many tasks which can be executed in parallel. The main components involved are: 1) a set of N compute nodes, each with S slots (one slot can run one task at a time), 2) a set of J jobs, each with T tasks, and 3) a set of D blocks of data, distributed across the compute nodes, where each task in a job operates on one (specific) block.

To constitute fairness, each of the J jobs must have an equal share of the cluster’s resources. If there are 2 jobs, then each gets half of the cluster, if there are 3, then a third, etc. To implement this, a fair scheduler must maintain a sorted queue of the remaining jobs, based on their current share of the cluster. The job with the least running tasks (the job farthest below its fair share of resources) is at the head of this queue. Let us call this job $J_{priority}$. Tasks from this job will take priority over the rest in the queue.

An observation that can be made is that, when a compute slot opens up on one of the N nodes, there might not be any data blocks on the node that the tasks of $J_{priority}$ need. This is called *data locality*. A task can be considered to be data local when it is scheduled in a slot on a node collocated with its data. Fairness enforcement compromises data locality, because the next job in the scheduler’s queue (based on fairness priority) might not have data in the next open slot in the cluster. Tasks that run non-locally must pull their input data in from across the network, causing an increase in completion time.

Traditional fixed delay scheduling operates in the following way: if $J_{priority}$ does not have any data blocks it needs in the slot currently being considered, it will be temporarily skipped, and a task from a job further down the queue (a lower priority job) will be scheduled instead. $J_{priority}$ will remain at the head of the queue (because it is still farthest below its fair share). $J_{priority}$ will thus be the highest priority job for the next slot that opens up, which hopefully will have data it needs. After a finite number of failed scheduling

attempts (skips), $J_priority$ will be scheduled in the next open slot regardless of locality. The maximum number of skips is usually determined by the number of scheduling attempts that occur within a certain interval of time (a few seconds). We can call this interval of time the *delay interval*.

The argument behind why fixed delay scheduling is acceptable is a probabilistic one [8]. Assuming that slots free up frequently, the chance that an open slot will not overlap with a job’s input data decreases exponentially. This chance is further decreased by block replication in distributed file systems (like HDFS), each of which can serve data for a task; however, this probabilistic approach relies on both a) an equal probability of slots opening up on each machine across the cluster, and b) a constant slot opening frequency. If the frequency were to change (for example, if network conditions were to worsen and currently running tasks were to take longer to complete), then the chance of achieving locality during the same fixed interval decreases. As well, variance in job profiles can lead to skewed probabilities of slot openings from machine to machine (some tasks take longer than others, some jobs have data on only a subset of machines, etc.).

3. FIXED DELAY SCHEDULING IN CLUSTER FRAMEWORKS

In real systems, fixed delay scheduling is usually implemented using a preset delay interval (as a unit of time, in milliseconds), and timestamps, the first of which is taken the first time a job is considered for scheduling (whether a task is actually scheduled or not). This timestamp is checked whenever a slot opens. If the slot is on a machine without input data for a job, then tasks from that job may only be scheduled there if the time between the current time and the timestamp is greater than whatever the fixed delay interval is. If the slot *does* have input data for the job, then tasks can be scheduled no matter the timestamp. Regardless of how/why a task was scheduled for a job, once a task is scheduled then the timestamp is reset, and the process begins again. Thus, this approach requires very little extra state: one additional field for each job that is currently in progress. For the purposes of this paper, we will be focusing on how fixed delay scheduling operates in Apache Spark (detailed below), but the principles and implementations are the same in other systems, such as Hadoop, YARN, and Mesos.

Apache Spark [9] is a general execution engine for distributed “Big Data” processing, similar to Hadoop, that breaks jobs into stages of tasks, and uses fixed delay scheduling to improve data locality when launching these tasks. Spark provides many advantages, such as an abstraction for representing distributed datasets (Resilient Distributed Datasets, or RDDs). These RDDs, which represent data either on disk (in Hadoop Distributed File System, or HDFS) or in memory, are divided into partitions across worker processes called executors in the cluster. Each executor has a finite number of slots (usually equal to the number of CPU cores available) with which to run tasks. The Spark scheduler, which resides in a centralized driver process, uses fixed delay scheduling to provide greater data locality when scheduling tasks into slots that become open. The fixed delay interval (the maximum waiting time) is set statically in a configuration file.

Spark (and the other systems mentioned above) also pro-

vide hierarchical fixed delay scheduling. A delay interval can be specified for node-local, rack-local, or even datacenter-local scheduling. If the delay interval expires for a job, it will fall back into the next tier (e.g. from node to rack local). The job will delay there (for perhaps a different period of time than the previous interval) for slots with the relaxed locality condition of that tier (slots within the same rack, for example).

4. DYNAMIC DELAY SCHEDULING

4.1 Beyond Constant Delays

The frameworks described above use a fixed delay interval, which remains the same throughout an entire job. Our observation is that there is a calculable tipping point after which waiting/skipping is no longer beneficial. Generally, we can define the total running time of a task t as:

$$t_{total} = t_{compute} + t_{netOH} + sch_delay$$

where $t_{compute}$ is the local computation time of the task, t_{netOH} is the network latency incurred from data transfer (this is zero if the task is scheduled locally), and sch_delay is the amount of time t is delayed (skipped) to wait for locality, from the moment t is first considered for scheduling.

Ideally, we want to minimize the task completion time of all tasks within a given job, while also providing some guarantee as to the maximum latency introduced by the scheduler itself through delay scheduling. The goal of delay scheduling is to remove t_{netOH} by scheduling a task locally, but if the task is delayed longer than this t_{netOH} , then it would have actually been more efficient in hindsight (in terms of task completion time) to immediately schedule the task in the first open slot that was considered. Conversely, if the delay interval is too short, then the task could potentially miss scheduling opportunities with locality, which would result in a shorter completion time than scheduling without locality.

Dynamic Delay Scheduling improves upon fixed delay scheduling by dynamically adapting the maximum sch_delay equal to t_{netOH} on a per-task basis. This allows for the longest waiting time for data-local slots while limiting the worst case task running time to $t_{compute} + (2 * t_{netOH})$. In this way, the delay interval will adapt to the severity of the network overhead. If t_{netOH} is low, then locality is less of a concern and the scheduler can shorten the time it delays. If t_{netOH} becomes high, then the scheduler increases the time that it waits, as locality is more important in poorer network conditions; however, this means that the network overhead incurred by running tasks needs to be monitored and reported to the scheduler over time, in order to adapt to changing network conditions and potentially heterogeneous data block sizes.

4.2 An Adaptive Solution Using Task Feedback

The network overhead incurred by missing data locality depends on many factors, including network traffic, the distance of a task from its data, and the input data size. These factors vary from job to job, and some (like network traffic) can change during the job execution itself. In order to adapt the delay scheduling interval to changing conditions, we have designed a simple feedback mechanism, in which each task reports the network overhead it experienced upon

completion (t_{netOH} , for a task t). This metric is sent to a centralized fair scheduler, which then changes its delay scheduling interval using a running average of feedback from completed tasks. This new delay scheduling interval is then used when scheduling subsequent tasks of the same type.

The feedback metric itself (t_{netOH}) can be any time measurement which reflects the network overhead incurred by a task. In the case of Dynamic Delay Scheduling, we are considering tasks which read data from distributed file systems, so the t_{netOH} for each task is the amount of time spent reading data blocks remotely; however, any arbitrary type of task, which has some way of measuring its own overhead, could set this as well and benefit from the adaptation.

Algorithm 1: Dynamic Delay Scheduling

```

1 Initialization:
2 delay_interval ← default (3 seconds)
3
4 Upon receiving task t completion event:
5 delay_interval ← (delay_interval + t_netOH) / 2
6
7 Scheduling Procedure:
8 for each open slot s on a node n in N do
9   sort J in increasing order of currently running tasks
10  for j in J do
11    if j has just launched a task then
12      j.wait_begin = current_time
13    end
14    if a task t in j has data on s's node n then
15      launch t in s
16    else
17      if a task t in j is unlaunched then
18        if current_time - j.wait_begin >
19          delay_interval then
20          launch t in slot s
21        else
22          continue
23        end
24      end
25    end
26 end

```

Algorithm 1 describes our adaptive solution in detail, in a cluster with N nodes and J jobs. The algorithm consists of two parts: firstly, when tasks complete, the delay interval will be changed to reflect recent overhead measurements (line 5 of Algorithm 1). The delay interval adapts relatively quickly (the average of the last few readings), since feedback can only be received when tasks complete. For workloads in which tasks complete very often, the window for the average can be broadened. Secondly, this new delay interval is used when scheduling the next task (lines 7-26). The scheduler will choose the highest priority job, and check to see if it has data on the machine that is free. If it does, then we can launch a task locally. If it does not, then a check is done to determine if we are allowed to launch a task without locality. If the job has been waiting longer than the current delay interval (line 18), then we can launch a task. Otherwise, we skip this job (i.e. continue waiting) and continue on to the next job. Regarding fairness, each job only delays as long as the network overhead that would be experienced normally, and the decrease in job latency makes any temporary fairness imbalance insubstantial. Since a delayed job stays at the head of the scheduling queue while waiting, it will imme-

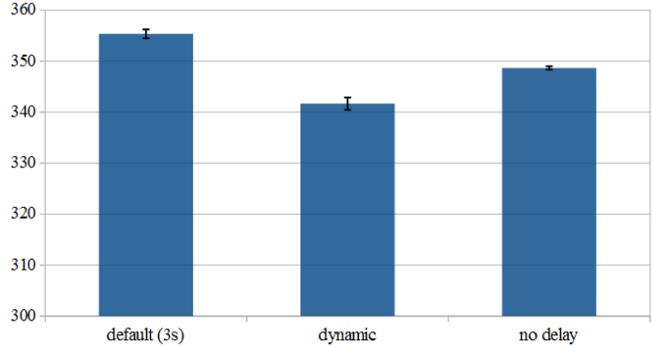


Figure 2: 1GB Spark Sort - Job Completion Time (s) vs. Delay Type

diately be scheduled up to its fair share if the delay interval expires.

It is worth noting that there is nothing preventing this adaptive solution from also being hierarchical, as many frameworks desire/implement separate delay intervals for node-local or rack-local scheduling.

5. IMPLEMENTATION IN SPARK

We implemented Dynamic Delay Scheduling in Spark, for tasks which read data from HDFS as input. Our implementation consists of about 100 lines of Scala code, which modify the already-existing fixed delay scheduling mechanisms as well as store a small bit of extra state.

The t_{netOH} for each task is implemented as a new field in each TaskContext, called $netOH$. To compute the overhead, we have inserted a wrapper for the HDFS RecordReader, which sums reading time. When a task completes, it will set $netOH$ to this sum, and then send it back to the Spark scheduler via the event bus that drives the system. While we have focused our implementation on “map-like” tasks which read from HDFS (which have predictable behavior), any arbitrary task could monitor its own overhead and set this field before closing.

The Spark Scheduler, when it sees a TaskCompleted event, will pull off that task’s $netOH$ and average it into a global $delay_interval$ that is used for all tasks in the same “stage” (stages are groups of tasks which do the same thing). Then, when a slot opens up in the cluster, the most recent $delay_interval$ will be used in place of the delay interval fetched from the spark-defaults.conf file.

Finally, we note that our implementation is simple enough to be adapted to other frameworks that use delay scheduling as well, as most systems (such as Hadoop) also have events/messages that inform the scheduler with information about finished tasks.

6. EVALUATION

In order to test the efficiency of Dynamic Delay Scheduling, we have designed several tests, on both a small and a large scale, using our implementation in Spark, with fair scheduling enabled. The goal of these experiments is to demonstrate that an adaptive delay interval provides decreased job latency in a variety of cluster setups and workloads.

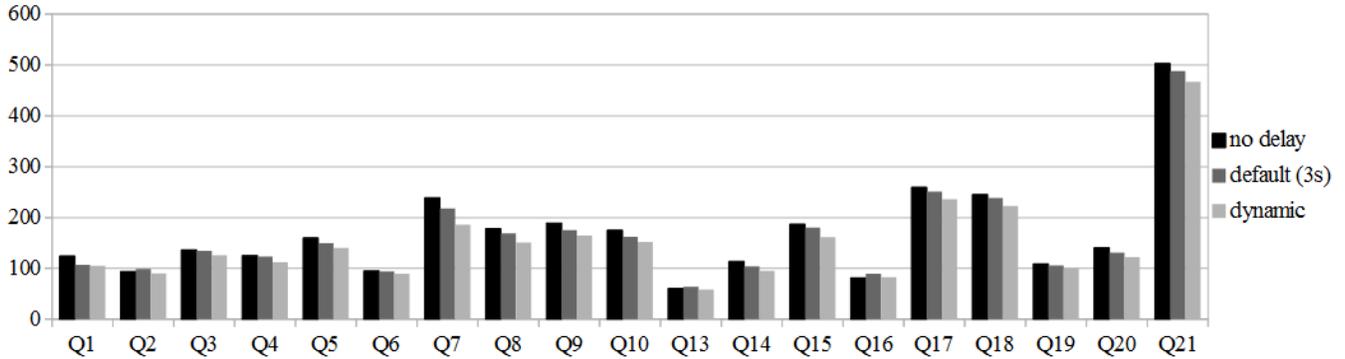


Figure 3: 110GB Spark TPC-H Queries - Job Completion Time (s) vs. Delay Type

6.1 Evaluation Setup

Our small-scale tests were performed on a four node local cluster, with four task slots (four 1.6 GHz cores for Spark tasks) per machine. All four nodes were connected through a Gigabit Ethernet switch. The workload chosen for the small-scale test was an integer sort, on 1GB of input data spread across HDFS. The main goal behind the small-scale setup is to evaluate our algorithm in a controlled environment with stable network conditions. To provide data locality concerns, input data was placed on two nodes, so as to ensure both data-local and non-local slots for task placement.

In order to have a more realistic evaluation, our large-scale experiments were performed on a 16 node cluster on Amazon EC2, using m3.xlarge instances. The workload chosen for this experiment was a TPC-H benchmark, which performs queries reading from multiple databases in HDFS, of size ranging from several megabytes to 80 gigabytes (for a total of 110GB of input data). We chose this particular benchmark because each query consists of a mixture of large and small stages which read from and process these tables in parallel. This causes contention for resources in the cluster, creating situations in which locality and fairness clash. This benchmark also represents a workload in which fairness is likely to be desired. Without fair sharing, stages which read from small tables would be significantly delayed by longer-running stages operating on larger data sets. In our experiments, queries 11, 12, and 22 failed to terminate in any configuration (perhaps due to bugs in the query implementation), and have been omitted from the figure. All numbers, from both small and large scale, are the average of 3 runs.

6.2 Small Scale Results

The results of our small-scale experiments are shown in Figure 2. We chose sorting because sorting is an IO-bound job, meaning that the effects of delay scheduling (increasing data locality) are more important to overall job performance. We performed the sort first with no delay at all, then with the default fixed 3 second delay, and then with dynamic delay scheduling enabled. The default of 3 seconds causes a slight increase in completion time (versus no delay). Given that network conditions were very good between the nodes in the cluster (all were connected to the same switch), having a fixed delay value of 3 seconds causes increased overhead, because the network overhead of missing locality is far shorter than 3 seconds. Our adaptive solution

performs slightly better, because the delay interval shrinks based on the feedback from the first round of tasks, although the job is still too small for there to be the kind of resource contention that would result in larger speedups. This experiment simply demonstrates that an incorrect fixed delay interval can actually hurt the performance of a job, rather than improve it.

6.3 Large Scale Results

The results of our large-scale TPC-H tests are shown in Figure 3, with the same three delay setups (no delay, fixed 3 seconds, and dynamic). In every query tested, Dynamic Delay Scheduling had lower job latency than the default fixed delay, by roughly 10% per query. With these being large workloads with hundreds of tasks, there is a great deal of contention for slots amongst all running stages, creating both more accurate feedback over time and more opportunities for delay scheduling to benefit job latency. We note that the benefit is largest in the longer queries, where there are both large and small stages competing for slots. In the shorter queries, we can see that a dynamic delay at least does not add any overhead, as the default delay does in queries Q2, Q13, and Q16. This shows that Dynamic Delay Scheduling can lower job latency by both a) removing excess waiting time in small jobs, and by b) improving fairness vs. performance decisions in jobs with high resource contention.

7. RELATED WORK

In the domain of job scheduling in data analytics frameworks like Spark, delay scheduling [8] has been suggested as the reasonable strategy for achieving fair sharing for jobs and data locality for tasks. Some schedulers (like Hadoop YARN’s [6] Fair Scheduler) have begun implementing fixed delay scheduling as a skipcount proportional to the size (in number of nodes) of the compute cluster. This approach adapts the interval to the size of the cluster, but not to the usage or load of the cluster in real time. Additionally, YARN considers coarse-grained resource allocations for entire frameworks, as opposed to the finer-grained tasks we consider.

Sparrow [4] is a distributed task scheduler, which focuses on a particular set of high-frequency, interactive workloads. Sparrow features distributed, stateless scheduling and randomized load balancing of tasks to reduce job latency; however, its completely distributed nature makes enforcing fairness policies challenging, and its benefits decrease as work-

loads stray from the interactive scale. Our goal, on the other hand, is to provide lower job latency for a diversity of workloads, through simple adaptive techniques.

There is also an extensive amount of work on scheduling for custom cluster-management frameworks that attempt to minimize task completion times. Quincy [3] is a scheduling system which explores the tradeoffs of various types of fairness policy enforcement, in the presence of data locality concerns. In particular, Quincy focuses on the use of preemption, to enforce either fairness or locality based on a computed optimal schedule. One caveat to this approach is that computing/recalculating a complete schedule causes increased overhead to making scheduling decisions. Our solution seeks to improve scheduling with very little extra overhead or state, keeping decisions fast. Microsoft Apollo [1] is a scheduling framework which uses distributed, loosely synchronized schedulers to opportunistically schedule tasks, correcting conflicts should they arise. Apollo features a service which advertises load to schedulers to facilitate better scheduling decisions. Unlike our system, Apollo’s architecture has tight integration between resource management and job/task scheduling. We consider systems that schedule between/within jobs, which have to efficiently use resources already acquired and shared amongst running jobs. Both Quincy and Apollo require a large amount of extra state (services that monitor load in real time, solvers to compute/recompute schedules, etc.). We hope to introduce a lightweight adaptive solution that can be applied to existing systems without the need for strenuous additional state.

8. CONCLUDING REMARKS

Delay scheduling has long been accepted as a *de facto* solution for achieving higher data locality (and assumedly better performance) in cluster framework schedulers that enforce fairness. We show that delay scheduling, using a fixed delay interval, is inefficient, and can be improved by adapting the delay interval on-the-fly, based on the overhead experienced by completed tasks. An adaptive solution allows delay scheduling to find a proper balance between fairness and performance, based on changing network conditions, job profiles, and cluster usage. We present Dynamic Delay Scheduling as a simple proof of concept that adaptive scheduling in these systems can provide benefits with minimal state.

Regarding future work, there are other areas where adaptivity looks promising in the realm of delay scheduling. For example, the frequency with which tasks complete is major contributor to the probability that delay scheduling will actually provide locality, and could potentially be used as another adaptive improvement. Other feedback metrics (such as task completion time) hold promise for improvements/heuristics as well. We plan to evaluate these adaptive approaches with other varied workloads, such as TPC-DS. Another promising direction is the analysis of the applicability of this adaptivity to stream processing frameworks, such as Spark Streaming [10].

9. REFERENCES

- [1] Eric Boutin, Jaliya Ekanayake, Wei Lin, Bing Shi, Jingren Zhou, Zhengping Qian, Ming Wu, and Lidong Zhou. Apollo: Scalable and coordinated scheduling for cloud-scale computing. In *Symp. on Op. Sys. Design and Implementation (OSDI)*, pages 285–300, 2014.
- [2] Benjamin Hindman, Andy Konwinski, Matei Zaharia, Ali Ghodsi, Anthony D. Joseph, Randy Katz, Scott Shenker, and Ion Stoica. Mesos: A platform for fine-grained resource sharing in the data center. In *Symp. on Op. Sys. Design and Implementation (OSDI)*, page 32, 2011.
- [3] Michael Isard, Vijayan Prabhakaran, Jon Currey, Udi Wieder, Kunal Talwar, and Andrew Goldberg. Quincy: Fair scheduling for distributed computing clusters. In *Symp. on Op. Sys. Principles (SOSP)*, pages 261–276, 2009.
- [4] Kay Ousterhout, Patrick Wendell, Matei Zaharia, and Ion Stoica. Sparrow: Distributed, low latency scheduling. In *Symp. on Op. Sys. Principles (SOSP)*, pages 69–84, November 2013.
- [5] Konstantin Shvachko, Hairong Kuang, Sanjay Radia, and Robert Chansler. The hadoop distributed file system. In *Symp. on Mass Storage Systems and Technologies (MSST)*, 2010.
- [6] Vinod Kumar Vavilapalli, Arun C Murthy, Chris Douglas, Sharad Agarwal, Mahadev Konar, Robert Evans, Thomas Graves, Jason Lowe, Hitesh Shah, Siddharth Seth, Bikas Saha, Carlo Curino, Owen O’Malley, Sanjay Radia, Benjamin Reed, and Eric Baldeschwieler. Apache hadoop yarn : Yet another resource negotiator. In *Symp. on Cloud Computing (SoCC)*, page 16, 2013.
- [7] Matei Zaharia, D Borthakur, J S Sarma, K Elmeleegy, S Shenker, and I Stoica. Job scheduling for multi-user mapreduce clusters. *EECS Department University of California Berkeley Tech Rep*, (UCB/EECS-2009-55):2009–55, 2009.
- [8] Matei Zaharia, Dhruba Borthakur, Joydeep Sen Sarma, Khaled Elmeleegy, Scott Shenker, and Ion Stoica. Delay scheduling: a simple technique for achieving locality and fairness in cluster scheduling. In *Euro. Conf. on Comp. Sys. (EuroSys)*, pages 265–278, 2010.
- [9] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, and Ankur Dave. Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing. In *Networked Sys. Design and Implem. (NSDI)*, pages 15–28, 2012.
- [10] Matei Zaharia, Tathagata Das, Haoyuan Li, Timothy Hunter, Scott Shenker, and Ion Stoica. Discretized streams: Fault-tolerant streaming computation at scale. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 423–438, New York, NY, USA, 2013. ACM.
- [11] Matei Zaharia, Andy Konwinski, Anthony D Joseph, Randy Katz, and Ion Stoica. Improving mapreduce performance in heterogeneous environments. In *Symp. on Op. Sys. Design and Implementation (OSDI)*, pages 29–42, 2008.