

Do the Hard Stuff First: Scheduling Dependent Computations in Data-Analytics Clusters

Robert Grandl[†], Srikanth Kandula, Sriram Rao, Aditya Akella[†], Janardhan Kulkarni
Microsoft and University of Wisconsin-Madison[†]
MSR-TR-2016-19

Abstract– We present D_{PS}, a scheduler that improves cluster utilization and job completion times by packing tasks with multi-resource requirements and inter-dependencies. While the underlying scheduling problem is intractable in general, D_{PS} is nearly optimal on the job DAGs that appear in production clusters at a large enterprise. Our key insight is that carefully handling the long-running tasks and those with tough-to-pack resource requirements will lead to good schedules for DAGs. However, which subset of tasks to treat carefully is *a priori* unclear. D_{PS} offers a novel search procedure that evaluates various possibilities and outputs a valid schedule. An online component enforces the schedules desired by the various jobs running on the cluster. In addition, it packs tasks and, for any desired fairness scheme, guarantees bounded unfairness. We evaluate D_{PS} on a server cluster using traces of over 100,000 DAGs collected from a large production cluster. Relative to the state-of-the-art schedulers, D_{PS} speeds up half of the jobs by over 50%.

1. INTRODUCTION

DAGs (directed acyclic graphs) are a powerfully general abstraction for scheduling problems. Scheduling network transfers of a multi-way join or the work in a geo-distributed analytics job and many others can be represented as DAGs. However, scheduling even one DAG is known to be an NP-hard problem [1, 2].

Consequently, existing work focuses on special cases of the DAG scheduling problem using simplifying assumptions such as: ignore dependencies, only consider chains, assume only two types of resources or only one machine or that the vertices have similar resource requirements [3, 4, 5, 6, 7, 8, 9]. However, the assumptions that underlie these approaches often do not hold in practical settings, motivating us to take a fresh look at this problem.

We illustrate the challenges in the context of job DAGs in data-analytics clusters. Here, each DAG vertex represents a computational task and edges encode input-output dependencies. Programming models such as SparkSQL, Dryad and Tez [10, 11, 12] lead to job DAGs that violate many of the above assumptions. Traces from a large cluster reveal that (a) DAGs have complex structures with the median job having a depth of seven and a thousand tasks, (b) there is substantial variation in compute, memory, network and disk usages across tasks (stdev./avg in requirements is nearly 100%), (c) task run-times range from sub-second to hundreds of seconds, and (d) clusters suffer from resource fragmentation across machines.

The net effect of these challenges, based on our analysis, is that the completion times of jobs in this production cluster can be improved by 50% for half the DAGs.

This problem is important because data-analytics clusters run thousands of mission critical jobs each day in enterprises and in the cloud [13, 14]. Even modest improvements in job throughput significantly improves the ROI (return-on-investment) of these clusters; and quicker job completion reduces the lag between data collection and decisions (i.e., “time to insight”) which potentially increases revenue [15].

To identify a good schedule for one DAG, we observe that the pathologically bad schedules in today’s approaches mostly arise due to these reasons: (a) long-running tasks have no other work to overlap with them and (b) the tasks that are runnable do not pack well with each other. Our core idea, in response, is rather simple: identify the potentially *troublesome* tasks, such as those that run for a very long time or are hard to pack, and place them first on a *virtual resource-time space*. This space would have $d + 1$ dimensions when tasks require d resources; the last dimension being time. Our claim is that placing the troublesome tasks first leads to a good schedule since the remaining tasks can be placed into resultant holes in this space.

Unfortunately, scheduling one DAG well does not suffice. Production cluster schedulers have many concurrent jobs, online arrivals and short-lived tasks [16, 17, 18, 19]. Together, these impose a strict time-budget on scheduling. Also, sharing criteria such as fairness have to be addressed during scheduling. Hence, production clusters are forced to use simple, online heuristics.

We ask whether it is possible to efficiently schedule complex DAGs *while* retaining the advantageous properties of today’s production schedulers such as reacting in an online manner, considering multiple objectives etc.

To this end, we design a new cluster scheduler D_{PS}. At job submission time or soon thereafter, D_{PS} builds a preferred schedule for a single job DAG by placing the troublesome tasks first. D_{PS} solves two key challenges in realizing this idea: (1) the best choice of troublesome tasks is intractable to compute and (2) *dead-ends* may arise because tasks are placed out-of-order (e.g., troublesome go first) and it is *a priori* unclear how much slack space should be set aside. D_{PS} employs a performant search procedure to address the first challenge and has a placement procedure that provably avoids dead-ends for the second challenge. Figure 1 shows an

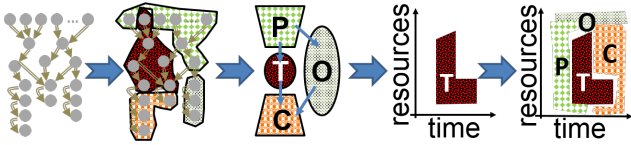


Figure 1: Shows steps taken by DAGPS from a DAG on the left to its schedule on the right. Troublesome tasks T (in red) are placed first. The remaining tasks (parents P , children C and other O) are placed *on top* of T in a careful order to ensure compactness and respect dependencies.

example.

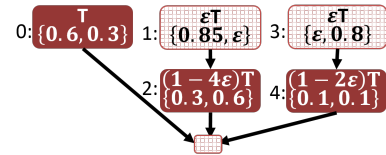
Each schedule constructed for each DAG is passed on to a second online component of DAGPS which coordinates between the various DAGs running in the cluster and also reconciles between their multiple, potentially discordant, objectives. For example, a fairness scheme such as DRF may require a certain job to get resources next, but multi-resource packing—which we use to reduce resource fragmentation—or the preferred schedules above may indicate that some other task should be picked next. Our reconciliation heuristic, colloquially, attempts to *follow the majority*; that is it can violate an objective, say fairness, when multiple other objectives counterweight it. However, to maintain predictable performance, our reconciliation heuristic limits maximum unfairness to an operator-configured threshold.

We have implemented the two components of DAGPS in Apache YARN and Tez and have experimented with jobs from TPC-DS, TPC-H and other benchmarks on a server cluster. Further, we also evaluate DAGPS in simulations on production DAGs from a production cluster.

To summarize, we make theoretical as well as practical contributions in this work. Our key contributions are:

- A characterization of the DAGs seen in production at a large enterprise and an analysis of the performance of various DAG scheduling algorithms (§ 2.1).
- A novel DAG scheduler that combines multi-resource packing and dependency awareness (§ 2.2).
- An online scheduler that mimics the preferred schedules for all the jobs on the cluster while bounding unfairness (§ 2.3) for many models of fairness [1, 2, 3].
- A new lower bound on the completion time of a DAG (§ 2.4). Using this we show that the schedules built by DAGPS’s online component are within 1.5 times OPT for half of the production DAGs; three quarters are within 2 times and the worst is 3 times OPT.
- An implementation that we intend to release as open source (§ 2.5).
- Our experiments show that DAGPS improves the completion time of half of the DAGs by 20%; the number varies across benchmarks. The improvement for production DAGs is at the high end of the range because these DAGs are more complex and have diverse resource demands.

Lastly, while our work is presented in the context of cluster scheduling, as noted above, similar DAG scheduling problems arise in other domains. We offer early results in Sec-



Technique	Execution Order	Time	Worst-case
OPT	$\{t, t\} \rightarrow \{t, t, t\} \rightarrow$	T	—
CPSched	$t \rightarrow t \rightarrow t \rightarrow t \rightarrow t \rightarrow$	T	$O(n) \times \text{OPT}$
Tetris	$t \rightarrow t \rightarrow t \rightarrow t \rightarrow t \rightarrow$	T	$O(d) \times \text{OPT}$

Figure 2: An example DAG where Tetris [37] and Critical Path Scheduling take \times longer than the optimal algo OPT. Here, DAGPS equals OPT. Details are in §2.2. Assume $\epsilon \rightarrow 0$.

tion from applying DAGPS to scheduling DAGs arising in distributed build systems [4, 5] and in request-response workflows [6, 7].

2. PRIMER ON SCHEDULING JOB DAGS

2.1 Problem definition

Let each job be represented as a directed acyclic graph $\mathcal{G} = \{V, E\}$. Each node in V is a task with demands for various resources. Edges in E encode precedence constraints between tasks. Many jobs can simultaneously run in a cluster.

A cluster is a group of servers organized as per some network topology.

DAGPS considers task demands along four resource dimensions (cores, memory, disk and network bandwidth). Depending on placement, tasks may need resources at more than one machine (e.g., if input is remote) or along network paths. The network bottlenecks are near the edges (at the source or destination servers and top-of-rack switches) in today’s datacenter topologies [8, 9, 10]. Some systems require users to specify the DAG \mathcal{G} explicitly [11, 12] whereas others use a query optimizer to generate \mathcal{G} [13]. Production schedulers already allow users to specify task demands (e.g., [14 core, 1 GB] is the default for tasks in Hadoop [15]). Note that such annotation tends to be incomplete (network and disk usage is not specifiable) and is in practice significantly overestimated since tasks that exceed their specified usage will be killed. Similar to other reports [16, 17], up to 50% of the jobs in the examined cluster are *recurring*. For such jobs, DAGPS uses past job history to estimate task runtimes and resource needs. For the remaining ad-hoc jobs, DAGPS uses profiles from similar jobs and adapts these profiles online (see § 2.2).

Given a set of concurrent jobs $\{\mathcal{G}\}$, the cluster scheduler maps tasks on to machines while meeting resource capacity limits and dependencies between tasks. Improving **performance**—measured in terms of the job throughput (or makespan) and the average job completion time—is crucial, while also maintaining **fairness**—measured in terms of how resources are divided amongst groups of jobs per some requirement (e.g., DRF or slot-fairness).

2.2 An illustrative example

We use the DAG shown in Figure 2 to illustrate the scheduling issues. Each node represents one task: the node

labels represent the task duration (top) and the demands for two resources (bottom). Assume that the total resource available is ϵ for both resources and let ϵ represent a small value.

Intuitively, a good schedule would overlap the long-running tasks shown with a dark background. The resulting optimal schedule (OPT) is shown in the table (see Figure 3). OPT overlaps the execution of all the long-running tasks— t_1 , t_2 and t_3 —and finishes in T . However, such long-running/resource intensive tasks can be present anywhere in the DAG, and it is unlikely that greedy local schedulers can overlap these tasks. To compare, the table also shows the schedules generated by a typical DAG scheduler, and a state-of-the-art *packer* which carefully packs tasks onto machines to maximize resource utilization. We discuss them next.

DAG schedulers such as critical path scheduling (CPSched) pick tasks along the critical path (CP) in the DAG. The CP for a task is the longest path from the task to the job output. The figure also shows the task execution order with CPSched. CPSched ignores the resources needed by tasks and does not pack. Consequently, for this example, CPSched performs poorly because it does not schedule tasks that are not on the critical path first (such as t_4 , t_5) even though doing so reduces resource fragmentation by overlapping the long-running tasks.

On the other hand, packers such as, Tetris [10], pack tasks to machines by matching along multiple resource dimensions. Tetris greedily picks the task with the highest value of the dot product between task’s demand vector and the available resource vector. The figure also shows the task execution order with Tetris. Tetris does not account for dependencies. Its packing heuristic only considers the tasks that are currently schedulable. In this example, Tetris performs poorly because it will not choose locally inferior packing options (such as running t_4 instead of t_5) even when doing so can lead to a better global packing.

DPS achieves the optimal schedule for this example. When searching for troublesome subsets, it will consider the subset $\{t_1, t_2, t_3\}$ because these tasks run for much longer. As shown in Figure 3, the troublesome tasks will be placed first. Since there are no dependencies among them, they will run at the same time. The parents ($\{t_4, t_5\}$) and any children are then placed *on top*; i.e., compactly before and after the troublesome tasks.

2.3 Analyzing DAGs in Production

We examined the production jobs from a cluster of tens of thousands of servers at a large enterprise. We also analyzed jobs from a server cluster that ran Hive [11] jobs and jobs from a high performance computing cluster [12].

To quantify potential gains, we compare the runtime of the DAGs in production to three measures. The first measure, CPLength is the duration of the DAG’s critical path. If the

CP of t_1, t_2, t_3 is $T, T(-\epsilon)$ and $T(-\epsilon)$ respectively. The demands of these tasks ensure that they cannot run simultaneously.

Tetris’ packing score for each task, in descending order, is $t_1 = . . . , t_2 = . . . , t_3 = . . .$ and $t_4 = . . .$

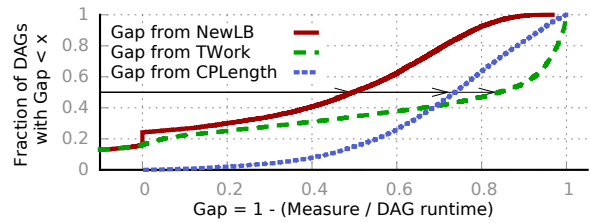


Figure 3: CDF of gap between DAG runtime and several measures. Gap is computed as $\frac{\text{measure}}{\text{DAG runtime}}$.

	CPU	Mem.	Network		Disk	
			Read	Write	Read	Write
Enterprise: Private Stack
Enterprise: Hive
HPC: Conductor	.	.	N/A	N/A	.	(R+W)

Table 1: Coefficient-of-variation (= stdev./avg.) of tasks’ demands for various resource. Across three examined frameworks, tasks exhibit substantial variability (CoV ~ .5) for many resources.

available parallelism is infinite, the DAG would finish within CPLength. The second measure, TWork, is the total work in the DAG normalized by the cluster share of that DAG. If there were no dependencies and perfect packing, a DAG would finish within TWork. In practice, both of these measures are quite loose—the first ignores all the work off the critical path and the second ignores dependencies. Hence, our third measure is a new improved lower bound NewLB that uses the specific structure of data-parallel DAGs. Further details are in §4 but intuitively NewLB leverages the fact that unlike random DAGs, all the tasks in a job stage (e.g., a map or reduce or join) have similar dependencies, durations and resource needs.

Figure 3 plots a CDF of the gap over all DAGs for these three measures. Observe that half of the jobs have a gap of over 0.5 for both CPLength and TWork. The gap relative to NewLB is smaller, indicating that the newer bound is tighter, but the gap is still over 0.5 for half of the jobs. That is, they take over two times longer than they could.

A few issues are worth noting for this result. First, some DAGs finish faster than their TWork and NewLB measures.

This is because our production scheduler is work conserving and can give jobs more than their fair share. Second, we know that jobs take longer in production because of runtime artifacts such as task failures or stragglers [13, 14]. What fraction of the gap is explained due to these reasons? When computing the job completion times to use in this result, we attempted to explicitly avoid these issues as follows. First, we chose the fastest completion time from among groups of related recurring jobs. It is unlikely that every execution suffers from failures. Second, we shorten the completion time of a job by deducting all periods when the job has fewer than ϵ tasks running concurrently. This explicitly corrects for stragglers—one or a few tasks holding up job progress. Hence, we believe that the remaining gap is likely due to the scheduler’s inability to pack tasks with dependencies.

To understand the causes for the performance gap further, we characterize the DAGs along the following dimensions:

"Work" that is ...	Percentage of total work in the DAG				
	[-)	[-)	[-)	[-)	[-)
on CriticalPath	41.6
"unconstrained"	33.3
"unordered"	56.6

Table 2: Bucketed histogram of where the work lies in DAGs. Each entry denotes the fraction of all DAGs that have the metric labeled on the row in the range denoted by the column. For example, 41.6% of DAGs have [-) of their total work on the critical path.

What do the DAGs look like? By depth, we refer to the number of tasks on the critical path. A map-reduce job has depth d . We find that the median DAG has depth $d/2$. Further, we find that the median ($d/2$ th percentile) task in-degree and out-degree are $(d/2)$ and $(d/2)$ respectively. If DAGs are chains of tasks, in- and out-degrees will be d . A more detailed characterization of DAGs including tree widths and path widths has been omitted for brevity. Our summary is that the vast majority of DAGs have complex structures.

How diverse are the resource demands of tasks? Table 2 shows the coefficient-of-variation (CoV) across tasks for various resources. We find that the resource demands vary substantially. The variability is possibly due to differences in work at each task: some are compute heavy (e.g., user-defined code that processes videos) whereas other tasks are memory heavy (e.g., in-memory sorts).

Where does the work lie in a DAG? We now focus on the more important parts of each DAG— the tasks that do more work (measured as the product of task duration and resource needs). Let CP_{Work} be the total work in the tasks that lie on the critical path. From Table 2, 41.6% of DAGs have CP_{Work} above 0.5. DAG-aware schedulers may do well for such DAGs. Let $UnconstrainedWork$ be the total work in tasks with no parents (i.e., no dependencies). We see that roughly 33.3% of the DAGs have $UnconstrainedWork$ above 0.5. Such DAGs will benefit from packers. The above cases are not mutually exclusive and together account for 74.9% of DAGs. For the other 25.1% of DAGs, neither packers nor criticality-based schedulers may work well.

Let $MaxUnorderedWork$ be the largest work in a set of tasks that are neither parents nor children of each other. Table 2 shows that 56.6% of DAGs have $MaxUnorderedWork$ above 0.5. That is, if ancestors of the unordered tasks were scheduled appropriately, substantial gains can accrue from packing the tasks in the maximal unordered subset.

From the above analysis, we observe that (1) production jobs have large DAGs that are neither a bunch of unrelated stages nor a chain of stages, and (2) a packing+dependency-aware scheduler can offer substantial improvements.

2.4 Analytical Results

We take a step back to offer some more general comments. First, DAG schedulers have to be aware of dependencies. That is, considering just the runnable tasks does not suffice.

Lemma 1. Any scheduling algorithm, deterministic or randomized, that does not account for the DAG structure is at least (d) times OPT where d is the number of resources.

For deterministic algorithms, the proof follows from de-

signing an adversarial DAG for any scheduler. We extend to randomized algorithms by using Yao's max-min principle (see A). Lemma 1 applies to all multi-resource packers [1, 2, 3, 4] since they ignore dependencies.

Second, and less formally, we note that schedulers have to be aware of resource heterogeneity. Many known scheduling algorithms have poor worst-case performance. In particular:

Lemma 2. Critical path scheduling can be (n) times OPT where n is the number of tasks in a DAG and Tetris can be $(d - 1)$ times OPT.

The proof is by designing adversarial DAGs for each scheduler (see B).

To place these results in context, note that d is about 10 (cores, memory, network, disk) and can be larger when tasks require resources at other servers or on many network links. Further, the median DAG has hundreds of tasks (n). D-PS is close to OPT on all of the described examples. Furthermore, D-PS is within 2 times optimal for half of the production DAGs (estimated using our new lower bound).

Finally, we note the following:

Lemma 3. If there were no precedence constraints and tasks were malleable, OPT is achievable by a greedy algorithm.

We say a task is malleable if assigning any (non-negative) portion of its demand p will cause it to make progress at rate p . In particular, tasks can be paused ($p = 0$) at any time which is also referred to as tasks being preemptible. The proof follows by describing the simple greedy algorithm which we omit here for brevity.

Our summary is that practical DAGs are hard to schedule because of their complex structure as well as discretization issues when tasks need multiple resources (fragmentation, task placement etc.)

3. NOVEL IDEAS IN D-PS

Cluster scheduling is the problem of matching tasks to machines. Every practical scheduler today does so in an online manner but has very tight timing constraints since clusters have thousands of servers, many jobs each having many pending tasks and tasks that finish in seconds or less [5, 6]. Given such stringent time budget, carefully considering large DAGs seems hopeless.

As noted in § 2, a key design decision in D-PS is to divide this problem into two parts. An offline component constructs careful schedules for a single DAG. We call these the *preferred schedules*. A second online component enforces the preferred schedules of the various jobs running in the cluster. We elaborate on each of these parts below. Figure 1 shows an example of how the two parts may inter-operate in a YARN-style architecture. Dividing a complex problem into parts and independently solving each part often leads to a sub-optimal solution. Unfortunately, we have no guarantees for our particular division. However, it can scale to large clusters and outperforms the state-of-art in experiments.

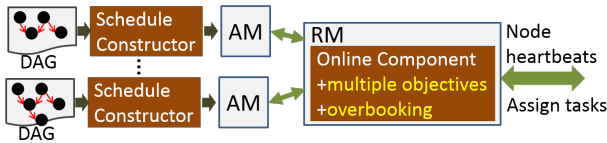


Figure 4: DAGPS builds schedules per DAG at job submission. The runtime component handles online aspects. AM and RM refer to the YARN’s application and resource manager components.

To find a compact schedule for a single DAG, our idea is to place the troublesome tasks, i.e., those that can lead to a poor schedule, first onto a virtual space. Intuitively, this maximizes the likelihood that any *holes*, un-used parts of the resource-time space, can be filled by other tasks. However, finding the best choice of troublesome tasks is as hard as finding a good schedule for the DAG. We use an efficient search strategy that mimics dynamic programming: it picks subsets that are more likely to be useful and avoids redundant exploration. Further, placing troublesome tasks first can lead to *dead-ends*. We define dead-end to be an arrangement of a subset of the DAG in the virtual space on which the remaining tasks cannot be placed without violating dependencies. Our strategy is to divide the DAG into subsets of tasks and place one subset at a time. While intra-subset dependencies are trivially handled by schedule construction, inter-subset dependencies are handled by restricting the order in which the various subsets are placed. We prove that the resultant placement has no *dead-ends*.

The online component has to co-ordinate between some potentially discordant directives. Each job running in the cluster offers a preferred schedule for its tasks (constructed as above). Fairness models such as DRF may dictate which job (or queue) should be served next. The set of tasks that is advantageous for packing (e.g., maximal use of multiple resources) can be different from both the above choices. We offer a simple method to reconcile these various directives. Our idea is to compute a real-valued score for each pending task that incorporates the above aspects *softly*. That is, the score trades-off violations on some directives if the other directives weigh strongly against it. For example, we can pick a task that is less useful from a packing perspective if it appears much earlier on the preferred schedule. Two key novel aspects are judiciously overbooking resources and bounding the extent of unfairness. Overbooking allows schedules that overload a machine or a network link if the cost of doing so (slowing-down of all tasks using that resource) is less than the benefit (can finish more tasks).

The online component of D_{PS} is described next; the online component is described in Section 4.

4. SCHEDULING ONE DAG

D_{PS} builds the schedule for a DAG in three steps. Figure 4 illustrates these steps and Figure 5 has a simplified pseudocode. First, D_{PS} identifies some troublesome tasks and divides the DAG into four subsets (§ 4.1). Second, tasks in a subset are packed greedily onto the virtual space while respecting dependencies (§ 4.2). Third, D_{PS} carefully restricts

Definitions: In DAG \mathcal{G} , t denotes a task and s denotes a stage, i.e., a collection of similar tasks.

Let \mathcal{V} denote all the stages (and hence the tasks) in \mathcal{G} .

Let $\mathcal{C}(s, \mathcal{G})$, $\mathcal{P}(s, \mathcal{G})$, $\mathcal{D}(s, \mathcal{G})$, $\mathcal{A}(s, \mathcal{G})$, $\mathcal{U}(s, \mathcal{G})$ denote the children, parents, descendants, ancestors and unordered neighbors of s in \mathcal{G} .

For clarity, $\mathcal{U}(s, \mathcal{G}) = \mathcal{V} - \mathcal{A}(s, \mathcal{G}) - \mathcal{D}(s, \mathcal{G}) - \{s\}$

```

1 Func: BuildSchedule:
2 Input:  $\mathcal{G}$ : a DAG,  $m$ : number of machines
3 Output: An ordered list of tasks  $t \in \mathcal{G}$ 
4  $S_{best} \leftarrow \emptyset$  // best schedule for  $\mathcal{G}$  thus far
5 foreach sets  $\{T, O, P, C\} \in \text{CandidateTroublesomeTasks}(\mathcal{G})$  do
6    $\text{Space } \mathcal{S} \leftarrow \text{CreateSpace}(m)$  //resource-time space
7    $\mathcal{S} \leftarrow \text{PlaceTasks}(T, \mathcal{S}, \mathcal{G})$  // trouble goes first
8    $\mathcal{S} \leftarrow \text{TrySubsetOrders}(\{\text{OCP, OPC, COP, POC}\}, \mathcal{S}, \mathcal{G})$ 
9   if  $\mathcal{S} < S_{best}$  then  $S_{best} \leftarrow \mathcal{S}$  //keep the best schedule;
10 return OrderTasks( $\mathcal{G}, S_{best}$ )

```

Figure 5: Pseudocode for constructing the schedule for a DAG. Helper methods are in Figure 6.

the order in which different subsets are placed such that the troublesome tasks go first and there are no dead-ends (§ 4.1). D_{PS} picks the most compact schedule after iterating over many choices for troublesome tasks. We discuss some enhancements in § 4.2. The resulting schedule is passed on to the online component (§ 4.2).

4.1 Searching for troublesome tasks

To identify *troublesome* tasks, D_{PS} computes two scores per task. The first, *LongScore*, divides the task duration by the maximum value across all tasks. Tasks with a higher score are more likely to be on the critical path and can benefit from being placed first because other work can overlap with them.

The second, *FragScore*, reflects the packability of tasks in a stage (e.g., a map or a reduce). It is computed by dividing the total work in a stage (T_{Work} defined in § 4.1) by how long a greedy packer will take to schedule that stage. Tasks that are more difficult to pack would have a lower *FragScore*. Given thresholds l and f , D_{PS} picks tasks with $\text{LongScore} \geq l$ or $\text{FragScore} \leq f$. Intuitively, doing so biases towards selecting tasks that are more likely to hurt the schedule because they are too long or too difficult to pack. D_{PS} iterates over different values for the l and f thresholds to find a compact schedule.

To speed up this search, () rather than choose the threshold values arbitrarily, D_{PS} picks values that are discriminative, i.e., those that allow different subsets of tasks to be considered as troublesome and () D_{PS} remembers the set of troublesome tasks that were already explored (by previous settings of the thresholds) so that it will construct a schedule only once per unique troublesome set.

As shown in Figure 6, the set T is a closure over the chosen troublesome tasks. That is, T contains the troublesome tasks and all tasks that lie on a path in the DAG between two troublesome tasks. The parent and child subsets P , C consist of tasks that are not in T but have a descendant or ancestor in T respectively. The subset O consists of the remaining tasks.

4.2 Compactly placing tasks

Given a subset of tasks and a partially occupied space, how

See Definitions atop Fig. .

```

1 Func: CandidateTroublesomeTasks:
2 Input: DAG  $\mathcal{G}$ ; Output: list  $\mathcal{L}$  of sets  $T, O, P, C$ 
   // choose a candidate set of troublesome tasks; per choice, divide  $\mathcal{G}$ 
   into four sets
3  $\mathcal{L} \leftarrow \emptyset$ 
4  $\forall v \in \mathcal{G}, \text{LongScore}(v) \leftarrow v.\text{duration} / \max_{v' \in \mathcal{G}} v'.\text{duration}$ 
5  $\forall v \in \mathcal{G}, v$  in stage  $s, \text{FragScore}(v) \leftarrow$ 
    $\text{TWork}(s) / \text{ExecutionTime}(s)$ 
6 foreach  $l \in \delta, \delta, \dots$  do
7   foreach  $f \in \delta, \delta, \dots$  do
8      $T \leftarrow \{v \in \mathcal{G} | \text{LongScore}(v) \geq l \text{ or } \text{FragScore}(v) \leq f\}$ 
9      $T \leftarrow \text{Closure}(T)$ 
10    if  $T \in \mathcal{L}$  then continue // ignore duplicates;
11     $P \leftarrow \bigcup_{v \in T} \mathcal{A}(v, \mathcal{G}); C \leftarrow \bigcup_{v \in T} \mathcal{D}(v, \mathcal{G});$ 
12     $\mathcal{L} \leftarrow \mathcal{L} \cup \{T, \mathcal{V} - T - P - C, P, C\}$ 

```

Figure 6: Identifying various candidates for troublesome tasks and dividing the DAG into four subsets.

best to pack the tasks while respecting dependencies? One can choose to place the parents first or the children first. We call these the *forward* and *backward* placements respectively. More formally, the forward placement recursively picks a task all of whose ancestors have already been placed on the space and puts it at the earliest possible time after its latest finishing ancestor. The backward placement is analogously defined. Intuitively, both placements respect dependencies but can lead to very different schedules since greedy packing yields different results based on which tasks are placed first. Figure 6:PlaceTasksF shows one way to do this. Traversing the tasks in either placement has $n \log n$ complexity for a subset of n tasks and if there are m machines, placing tasks greedily has $n \log(mn)$ complexity.

4.3 Subset orders that guarantee feasibility

For each division of DAG into subsets T, O, P, C, D DPS considers these four orders: TOCP, TOPC, TPOC or TCOP. That is, in the TOCP order, it first places all tasks in T , then tasks in O , then tasks in C and finally all tasks in P . Intuitively, this helps because the troublesome subset T is always placed first. Further, we will shortly prove that these are the only orders beginning with T that will avoid *dead-ends*.

A subtle issue is worth discussing. Only one of the forwards or backwards placements (described above in § 4.2) are appropriate for some subsets of tasks. For example, tasks in P cannot be placed *forwards* since some descendants of these tasks may already have been placed (such as those in T). As we saw above, the forwards placement places a task after its last finishing ancestor but ignores descendants and can hence violate dependencies if used for P . Analogously, tasks in C cannot be placed *backwards*. Tasks in O can be placed in one or both placements, depending on the inter-subset order. Finally, since the tasks in T are placed onto an empty space they can be placed either forwards or backwards. Formally, this logic is encoded in Figure 6:TrySubsetOrders. We prove the following lemma.

Lemma 4. (Correctness) The method described in § 4.2–§ 4.3 satisfies all dependencies and is free of *dead-ends*. (Completeness) Further, the method explores every order that

```

1 Func: PlaceTasksF: // forward placement
2 Inputs:  $\mathcal{V}$ : subset of tasks to be placed,  $\mathcal{S}$ : space (partially filled),  $\mathcal{G}$ : a DAG
3 Output: a new space with tasks in  $\mathcal{V}$  placed atop  $\mathcal{S}$ 
4  $\mathcal{S} \leftarrow \text{Clone}(\mathcal{S})$ 
5 finished placement set  $\mathcal{F} \leftarrow \{v \in \mathcal{G} | v \text{ already placed in } \mathcal{S}\}$ 
6 while true do
7   ready set  $\mathcal{R} \leftarrow \{v \in \mathcal{V} - \mathcal{F} | \mathcal{P}(v, \mathcal{G}) \text{ already placed in } \mathcal{S}\}$ 
8   if  $\mathcal{R} = \emptyset$  then break // all done;
9    $v' \leftarrow$  task in  $\mathcal{R}$  with longest runtime
10   $t \leftarrow \max_{v \in \mathcal{P}(v', \mathcal{G})} \text{EndTime}(v, \mathcal{S})$ 
11  // place  $v'$  at earliest time  $\geq t$  when its resource needs can be met
12   $\mathcal{F} \leftarrow \mathcal{F} \cup v'$ 
13 Func: PlaceTasks( $\mathcal{V}, \mathcal{S}, \mathcal{G}$ ): // inputs and output are same as PlaceTasksF
14 return min(PlaceTasksF( $\mathcal{V}, \mathcal{S}, \mathcal{G}$ ), PlaceTasksB( $\mathcal{V}, \mathcal{S}, \mathcal{G}$ ))
15 Func: PlaceTasksB: // only backwards, analogous to PlaceTasksF.
16 Func: TrySubsetOrders:
17 Input:  $\mathcal{G}$ : a DAG,  $\mathcal{S}_{\text{in}}$ : space with tasks in  $T$  already placed
18 Output: Most compact placement of all tasks.
19  $\mathcal{S}, \mathcal{S}', \mathcal{S}'', \mathcal{S}''' \leftarrow \text{Clone}(\mathcal{S}_{\text{in}})$ 
20 return min( // pick the most compact among all feasible orders
21 PlaceTasksF( $\mathcal{C}, \text{PlaceTasksB}(\mathcal{P}, (\text{PlaceTasks}(\mathcal{O}, \mathcal{S}, \mathcal{G})), \mathcal{G}), \mathcal{G}$ ), // OPC
22 PlaceTasksB( $\mathcal{P}, \text{PlaceTasksF}(\mathcal{C}, (\text{PlaceTasks}(\mathcal{O}, \mathcal{S}, \mathcal{G})), \mathcal{G}), \mathcal{G}$ ), // OCP
23 PlaceTasksB( $\mathcal{P}, \text{PlaceTasksB}(\mathcal{O}, (\text{PlaceTasksF}(\mathcal{C}, \mathcal{S}, \mathcal{G})), \mathcal{G}), \mathcal{G}$ ), // COP
24 PlaceTasksF( $\mathcal{C}, \text{PlaceTasksF}(\mathcal{O}, (\text{PlaceTasksB}(\mathcal{P}, \mathcal{S}, \mathcal{G})), \mathcal{G}), \mathcal{G}$ ) // POC
25 );

```

Figure 7: Pseudocode for the functions described in § 4.2 and § 4.3.

places troublesome tasks first and is free of *dead-ends*.

We omit a detailed proof due to space constraints. Intuitively however, the proof follows from () all four subsets are closed and hence intra-subset dependencies are respected by both the placements in § 4.2, () the inter-subset orders and the corresponding restrictions to only use forwards and/or backwards placements specified in § 4.3 ensure dependencies across subsets are respected and finally, () every other order that begins with T either violates dependencies or leads to a dead-end (e.g., in TPCO, placing tasks in O can dead-end because some ancestors and descendants have already been placed).

4.4 Enhancements

We note a few enhancements. First, due to barriers it is possible to partition a DAG into parts that are totally ordered. Hence, any schedule for the DAG is a concatenation of per-partition schedules. This lowers complexity because one execution of BuildSchedule will be replaced by several executions each having fewer tasks. Some of the production DAGs can be split into four or more parts. Second, and along similar lines, whenever possible we reduce complexity by reasoning over stages. Stages are collections of tasks and are sometimes fewer in number than tasks. Finally, we carefully choose our data-structures (e.g., a time and resource indexed hash map of free regions in space) so that the most frequent operation, picking a region in resource-time space where a task will fit as described in § 4.3, can be executed efficiently.

5. SCHEDULING MANY DAGS

We describe our online algorithm that matches tasks to machines while co-ordinating discordant objectives: fairness, packing and enforcing the per-DAG schedules built by § 4. We offer the pseudocode in Figure 7 for complete-

```

1 Func: FindAppropriateTasksForMachine:
2 Input:  $\mathbf{m}$ : vector of available resources at machine;  $\mathcal{J}$ : set of jobs with task
  details  $\{t_{\text{duration}}, t_{\text{demands}}, t_{\text{priScore}}\}$ ;  $\text{deficit}$ : counters for fairness;
3 Parameters:  $\kappa$ : unfairness bound;  $\text{rp}$ : remote penalty
4 Output:  $\mathcal{S}$ , the set of tasks to be allocated on the machine
5  $\mathcal{S} \leftarrow \emptyset$ 
6 while true do
7   foreach task  $t$  do
8      $\{p\text{Score}_t, o\text{Score}_t\} \leftarrow \{ \cdot, \cdot \}$ 
9      $\text{rPenalty}_t \leftarrow t$  is locality sensitive?  $\text{rp}$ :
10    if  $t_{\text{demands}} \leq \mathbf{m}$  // fits? then
11       $p\text{Score}_t \leftarrow (\mathbf{m} \cdot t_{\text{demands}}) \text{rPenalty}_t$  // dot product
12    else
13       $\text{compute } o\text{Score}_t$  // overbooking score omitted for brevity.
14       $\text{job } j \ni t, \text{srpt}_j \leftarrow \sum \text{pending } u_{e_j} \cdot t_{\text{duration}} * |u_{\text{demands}}|$ 
15       $p\text{erfScore}_t \leftarrow t_{\text{priScore}} \{p\text{Score}_t, o\text{Score}_t\} - \eta \text{srpt}_j$ 
16     $t_{\text{best}} \leftarrow \arg \max \{p\text{erfScore}_t | t\}$  // task with highest perf score
17    if  $t_{\text{best}} = \emptyset$  then break // no new task can be scheduled on this
  machine;
18     $g' \leftarrow$  jobgroup with highest deficit counter
19    if  $\text{deficit}_{g'} \geq \kappa C$  then  $t_{\text{best}} \leftarrow \arg \max \{p\text{erfScore}_t | t \in g'\}$ ;
20     $\mathcal{S} \leftarrow \mathcal{S} \cup t_{\text{best}}$ 
21    // detail: reduce available resources  $\mathbf{m}$ .
22     $\text{deficit}_{g'} \leftarrow \text{deficit}_{g'} +$ 
     $f(t_{\text{demands}}) * \begin{cases} \text{fairShare}_{g'} - & t \in \text{jobgroup } g \\ \text{fairShare}_{g'} & \text{otherwise} \end{cases}$ 

```

Figure 8: Simplified pseudocode for the online component.

ness but focus only on () how the various objectives are coordinated and () how unfairness is bounded.

e pseudocode shows how various individual objectives are estimated. Packing score per task $p\text{Score}_t$ is a dot product between task demands and available resources []. Using remote resources un-necessarily, for example by scheduling a locality-sensitive task [] at another machine, is penalized by the value rPenalty_t . e value srpt_j estimates the remaining work in a job and is used to prefer short jobs which lowers average job completion time. We claim no novelty thus far. Suppose that t_{priScore} is the order over tasks required by the schedule from § ; t_{priScore} is computed by ranking tasks in increasing order of their begin time and then dividing the rank by the number of tasks in the DAG so that the value is between (task that begins first) and (for the last task).

An initial combination of the above goals happens in the computation of $p\text{erfScore}_t$. See the first box in Figure . A task will have non-zero $p\text{Score}_t$ only if its demands fit within available resources. Else, it can have a non-zero $o\text{Score}_t$ if it is worth overbooking. We use a lexicographic ordering between these two values. at is, tasks with non-zero $p\text{Score}$ beat any value of $o\text{Score}$. Multiplying with t_{priScore} steers the search towards tasks earlier in the constructed schedule. Finally, η is a parameter that is automatically updated based on the average srpt and $p\text{Score}$. Subtracting $\eta \cdot \text{srpt}_j$ prefers shorter jobs. Intuitively, the combined value $p\text{erfScore}_t$ so ly enforces the various objectives. For example, if some task is preferred by all individual objectives (belongs to shortest job, is most packable, is next in the preferred schedule), then it will have the highest $p\text{erfScore}$. When the objectives are discordant, colloquially, the task preferred

$$\text{CPLen}_G = \max_{\text{path } p \in \mathcal{G}} \sum_{\text{task } t \in p} t_{\text{duration}} \quad (a)$$

$$\text{TWork}_G = \max_{\text{resource } r} \frac{1}{C_r} \sum_{t \in \mathcal{G}} t_{\text{duration}} t_{\text{demands}} \quad (b)$$

$$\text{ModCP}_G = \max_{p \in \mathcal{G}} \max_{s \in p} (\max(\text{TWork}_s, \text{CPLen}_s) + \sum_{t \in p - \{s\}} \min_{t \in s'} t_{\text{dur}}.) \quad (c)$$

$$\text{NewLB}_G = \sum_{G' \in \text{Partitions}(G)} \max(\text{CPLen}_{G'}, \text{TWork}_{G'}, \text{ModCP}_{G'}) \quad (d)$$

Figure 9: Lower bound formulas for DAG \mathcal{G} ; p, s, t denote a path through the DAG, a stage and a task respectively. C , here, is the capacity available for this job. We developed ModCP and NewLB .

by a majority of objectives will have the highest $p\text{erfScore}$.

To bound unfairness, we use one additional step. We explicitly measure unfairness using deficit counters []. When the maximum unfairness (across jobgroups or queues) is above the specified threshold κC , where C is the cluster capacity, DPS picks only among tasks belonging to the most unfairly treated jobgroup. is is shown in the second box in Figure . Otherwise DPS picks the task with the highest $p\text{erfScore}$. It is easy to see that this bounds unfairness by κC . Further, we can support a variety of fairness schemes by choosing how to change the deficit counter. For example, choosing $f() =$ mimics slot fairness (see third box in Figure), and $f() =$ demand of the dominant resource mimics DRF [].

6. A NEW LOWER BOUND

We develop a new lower bound on the completion time of a DAG of tasks. As we saw in § . , previously known lower bounds are very loose. Since the optimal solution is intractable to compute, without a good lower bound, it is hard to assess the quality of a heuristic solution such as DPS .

Equations a and b describe the known bounds: critical path length CPLen and total work TWork . Equation d is (a simpler form of) our new lower bound. At a high level, the new lower bound uses some structural properties of these job DAGs. Recall that DAGs can be split into parts that are totally ordered (§ .). is lets us pick the best lower bound for each part independently. For a DAG that splits into a chain of tasks followed by a group of independent tasks, we could use CPLen of the chain plus the TWork of the group. A second idea is that on a path through the DAG, at least one stage has to complete entirely. at is, all of the tasks in some stage and at least one task in each other stage on the path have to complete entirely. is leads us to the ModCP_G formula in Equation c where one stage s along any path p is replaced with the total work in that stage. A few other ideas are omitted for brevity.

e take-away is that the new lower bound NewLB is much tighter and allows us to show that DPS is close to OPT ; since by definition of a lower bound $\text{DPS} \geq \text{OPT} \geq \text{NewLB}$.

7. DPS SYSTEM

We have implemented the runtime component (§) in the Apache YARN resource manager (RM) and the schedule constructor (§) in the Apache Tez application master (AM). Our

schedule constructor implementation finishes in tens of seconds on all of the DAGs used in experiments; this is in the same ballpark as the time to compile and query-optimize these DAGs. Further, recurring jobs use previously constructed schedules. Each DAG is managed by an instance of the Tez AM which closely resembles other popular frameworks such as FlumeJava [10] and Dryad [11]. The per-job AMs negotiate with the YARN RM for containers to run the job's tasks; each container is a fixed amount of various resources. As part of implementing DIPS, we expanded the interface between the AM and RM to pass additional information, such as the job's pending work and tasks' demands, duration and preferred order. Due to anonymity considerations, we are unable to share full details of our code release. Here, we describe two key implementation challenges: (a) constructing profiles of tasks' resource demands and duration (§ 7.1), and (b) efficiently implementing the new online task matching logic (§ 7.2).

7.1 Profiling Tasks' Requirements

We estimate and update the tasks' resource demands and durations as follows. Recurring jobs are fairly common in production clusters (up to [12, 13]), executing periodically on newly arriving data (e.g., updating metrics for a dashboard). For these jobs, DIPS extracts statistics from prior runs. In the absence of prior history, we rely on two aspects of data analytics computations that make it amenable to learn profiles at runtime. (1) Tasks in a stage (e.g., map or reduce) have similar profiles and (2) tasks often run in multiple waves due to capacity limits. DIPS measures the progress and resource usage of tasks at runtime. Using the measurements from in-progress and completed tasks, DIPS refines estimates for the remaining tasks. Our evaluation will demonstrate the effectiveness of this approach.

7.2 Efficient Online Matching: Bundling

We have redesigned the online scheduler in YARN that matches machines to tasks. From conversations with Hadoop committers, these code-changes help improve matching efficiency and code readability.

Some background: The matching logic is heartbeat based. When a machine heartbeats to the RM, the allocator (\mathcal{A}) picks an appropriate task to allocate to that machine, (\mathcal{A}) adjusts its data structures (such as, resorting/rescoring) and (\mathcal{A}) repeats these steps until all resources on the node have been allocated or all allocation requests have been satisfied.

As part of this work, we support *bundling* allocations. That is, rather than breaking the loop after finding the first schedulable task, we maintain a set of tasks that can all be potentially scheduled on the machine. This so-called *bundle* allows us to schedule multiple tasks in one iteration, admitting non-greedy choices over multiple tasks. For example, if tasks t_1, t_2, t_3 are discovered in that order, it may be better to schedule t_1 and t_2 together rather than schedule t_1 by itself. We refactored the scheduler to support bundling; with configurable choices for (\mathcal{A}) which tasks to add to the bundle, (\mathcal{A})

when to terminate bundling (e.g. the bundle has a good set of tasks) and (\mathcal{A}) which tasks to pick from the bundle.

8. EVALUATION

Here, we report results from experiments on a server cluster and extensive simulations using DAGs from production clusters. Our key findings are:

- (1) In experiments on a large server cluster, relative to Tez jobs running on YARN, DIPS improves completion time of half of the jobs by 20% to 40% across various benchmarks. A quarter of the jobs improve by 50% to 70%.
- (2) On the DAGs from production clusters, schedules constructed by DIPS are faster by 10% for half of the DAGs. A quarter of the DAGs improve by 20%. Further, by comparing with our new lower bound, these schedules are optimal for 30% of the jobs and within 10% of optimal for 60% of the jobs.

As part of the evaluation, we offer detailed comparisons with many alternative schedulers and sensitivity analysis to cluster load and parameter choices. We also provide early results on applying DIPS to DAGs from other domains (§ 8.2).

8.1 Setup

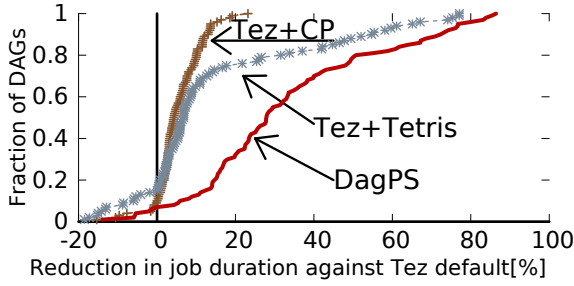
Our experimental cluster has 100 servers with two quad-core Intel E5 processors (hyperthreading enabled), 32 GB RAM, 8 drives, and a 10 Gbps network interface. The network has a congestion-free core [14].

Workload: Our workload mix consists of jobs from public benchmarks—TPC-H [15], TPC-DS [16], BigBench [17], and jobs from a production cluster that runs Hive jobs (E-Hive). We also use 1000 DAGs from a private production system in our simulations. In each experimental run, jobs arrival is modeled via a Poisson process with average inter-arrival time of 10s for 10 minutes. Each job is picked at random from the corresponding benchmark. We built representative inputs and varied input size from GBs to tens of TBs such that the average query completes in a few minutes and the longest finishes in under 10 minutes on the idle cluster. A typical experiment run thus has about 100 jobs and lasts until the last job finishes. The results presented are the median over three runs.

Compared Schemes: We experimentally compare DIPS against the following baselines: (1) Tez : breadth-first order of tasks in the DAG running atop YARN's Capacity Scheduler (CS), (2) Tez + CP : critical path length based order of tasks in the DAG atop CS and (3) Tez + Tetris : breadth-first order of tasks in the DAG atop Tetris [18].

Using simulations, we compare DIPS against the following schemes: (1) BFS : breadth-first order, (2) CP : critical path order, (3) Random order, (4) StripPart [19], (5) Tetris [18], and (6) Coffman – Graham [20].

All of the above schemes except (1) are work-conserving. (1)–(3) and (4) pick greedily from among the runnable tasks but vary in the specific heuristic. (5) and (6) require more complex schedule construction, as we will discuss later.



(a) CDF of gains for jobs on TPC-DS workload

Workload	th percentile			th percentile		
	D	T+C	T+T	D	T+C	T+T
TPC-DS	27.8	.	.	45.7	.	.
TPC-H	30.5	.	.	48.3	.	.
BigBench	25.0	.	.	33.3	.	.
E-Hive	19.0	.	.	29.7	.	.

D stands for D PS. T+C and T+T denote Tez + CP and Tez + Tetris respectively (see §). The improvements are relative to Tez.

(b) Improvements in job completion time across all the workloads

Figure 10: Comparing completion time improvements of various schemes relative to Tez.

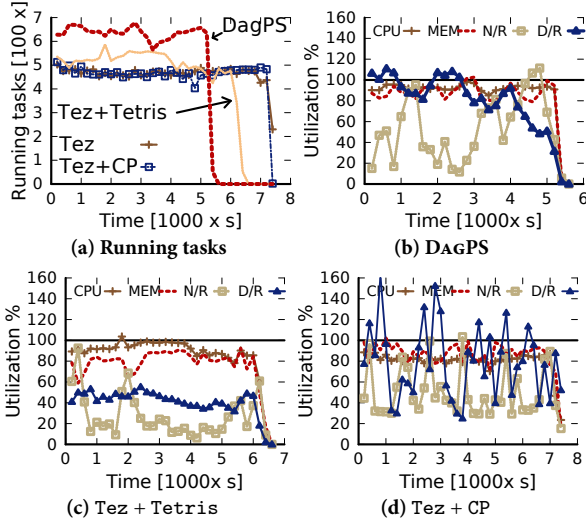


Figure 11: For a cluster run with jobs, a time lapse of how many tasks are running (leftmost) and how many resources are allocated by each scheme. N/R represents the amount of network read, D/R the disk read and D/W the corresponding disk write.

Metrics: Improvement in job completion time is our key metric. Between two schemes, we measure the *normalized gap* in job completion time. That is, the difference in the runtime achieved for the same job divided by the runtime of the job with some scheme; the normalization lets us compare across jobs with very different runtimes. Other metrics of interest are makespan, i.e., the time to finish a given set of jobs, and Jain’s fairness index [] to measure how close the cluster scheduler comes to the desired allocations.

8.2 How does D PS do in experiments?

8.2.1 Job Completion Time

Relative to Tez, Figure shows that D PS improves half of the DAGs by – across various benchmarks. One quarter of the DAGs improve by – . We see occasional

Workload	Tez+CP	Tez+Tetris	D PS
TPC-DS	+	+	+
TPC-H	+	+	+

Table 3: Makespan, gap from Tez.

Workload	Scheme	2Q vs. 1Q Perf. Gap	Jain’s fairness index		
			S	S	S
TPC-DS	Tez	–	.	.	.
	Tez+DRF	–	.	.	.
	Tez+Tetris	–	.	.	.
	DAGPS	+	.	.	.

Table 4: Fairness: Shows the performance gap and Jain’s fairness index when used with 2 queues (even share) versus 1 queue. Here, a fairness score of indicates perfect fairness.

regressions. Up to of the jobs slow down with D PS; the maximum slowdown is . We found this to be due to two reasons. (a) Noise from runtime artifacts such as stragglers and task failures and (b) Imprecise profiles: in all of our experiments, we use a single profile (the average) for all tasks in a stage but due to reasons such as data-skew, tasks in a stage can have different resource needs and durations. The table in Fig. shows results for other benchmarks; we see that DAGs from E-Hive see the smallest improvement (at median) because the DAGs here are mostly two stage map-reduce jobs. The other benchmarks have more complex DAGs and hence receive sizable gains.

Relative to the alternatives, Figure shows that D PS is to better. Tez + CP achieves only marginal gains over Tez, hinting that critical path scheduling does not suffice. The exception is the BigBench dataset where about half the queries are dominated by work on the critical path. Tez + Tetris comes closest to D PS because Tetris’ packing logic reduces fragmentation. But, the gap is still substantial, since Tetris ignores dependencies. In fact, we see that Tez + Tetris does not consistently beat Tez + CP. Our takeaway is that considering both dependencies and packing can substantially improve DAG completion time.

Where do the gains come from? Figure offers more detail on an example experimental run. D PS keeps more tasks running on the cluster and hence finishes faster (Fig. a).

The other schemes take over longer. To run more tasks, D PS gains by reducing fragmentation and by overbooking fungible resources. Comparing Fig. b with Figs. c– d, the average allocation of all resources is higher with D PS. Occasionally, D PS allocates over of the network and disk. Tez + Tetris, the closest alternative, has fewer tasks running at all times because (a) it does not overbook (all resource usages are below in Fig. c) and (b) it ignores dependencies and packs greedily leading to a worse packing of the entire DAG. Tez + CP is impacted negatively by two effects: (a) ignoring disk and network usage leads to arbitrary over-allocation (the “total” resource usage is higher because, due to saturation, tasks hold on to allocations for longer) and (b) due to fragmentation, many fewer tasks run on average. Together these lead to low task throughput and job delays.

8.2.2 Makespan

To evaluate makespan, we make one change to experiment setup— all jobs arrive within the first few minutes. Everything

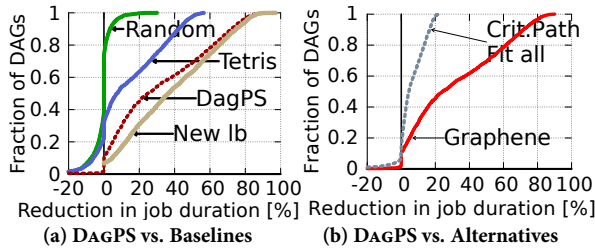


Figure 12: Comparing DAGPS with other schemes. We removed the lines from CG and StripPart because they hug $x = 0$; see Table 5.

else remains the same. Table 5 shows the gap in makespan for different cases. Due to careful packing, DAGPS sustains high cluster resource utilization, which in turn enables individual jobs to finish quickly: makespan improves relative to Tez and overall relative to alternatives.

8.2.3 Fairness

Can we improve performance while also being fair? Intuitively, fairness may hurt performance since the task scheduling order needed for high performance (e.g., packability or dependencies) differs from the order that ensures fairness. To evaluate fairness, we make one change to the experiment set up. The jobs are evenly and randomly distributed among two queues and the scheduler has to divide resources evenly.

Table 5 reports the gap in performance (median job completion time) for each scheme when run with two queues vs. one queue. We see that Tez, Tez + DRF and Tez + Tetris lose overall in performance relative to their one queue counterparts. The table shows that with two queues, DAGPS has a small gain (perhaps due to experimental noise). Hence, relatively, DAGPS performs even better than the alternatives if given more queues (the gap at one queue in Fig. 12a translates to a larger gap at two queues). But why? Table 5 also shows Jain’s fairness index computed over 1s, 5s and 10s time windows. We see that DAGPS is less fair at short timescales but is indistinguishable at larger time windows. This is because DAGPS is able to bound unfairness (§ 4.2); it leverages some short-term *slack* from precise fairness to make scheduling choices that improve performance.

8.3 Comparing with alternatives

We use simulations to compare a much wider set of best-of-breed algorithms (§ 4.1) on the much larger DAGs that ran in the production clusters. We mimic the actual dependencies, task durations and resource needs from the cluster.

Figure 13 compares the schedules constructed by DAGPS with that from other algorithms. Table 5 reads out the gaps at various percentiles. We observe that DAGPS’s gains at the end of schedule construction are about the same as those obtained at runtime (Figure 12). This is interesting because the runtime component only solely enforces the desired schedules from all the jobs running simultaneously in the cluster. It appears that any loss in performance from not adhering to the desired schedule are made up by the gains from better packing and trading off some short-term unfairness.

Second, DAGPS’s gains are considerable compared to the al-

		7 th	25 th	57 th	74 th
DAGPS		7	25	57	74
Random		–	–	–	–
Crit.Path	Fit cpu/mem	–	–	–	–
	Fit all	–	–	–	–
Tetris	Fit all	–	–	–	–
Strip Part.	Fit all	–	–	–	–
Co man-Graham.	Fit all	–	–	–	–
	Fit cpu/mem	–	–	–	–

Table 5: Reading out the gaps from Figure 12; comparing DAGPS vs. Alternatives. Each entry is the improvement relative to BFS.

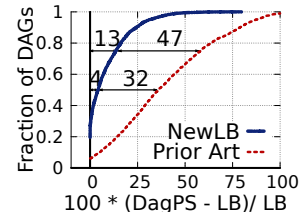


Figure 13: Comparing DAGPS with lower bounds.

ternatives. CP and Tetris are the closest. The reason is that DAGPS looks at the entire DAG and places the troublesome tasks first, leading to a more compact schedule overall.

Third, when tasks have unit durations and nicely shaped demands, CG (Co man-Graham [10]) is at most times optimal. However, it does not perform well on production DAGs that have diverse demands for resources and varying durations. Some recent extensions to CG handle heterogeneity but ignore fragmentation issues when resources are divided across many machines [11].

Fourth, StripPart [12] is the best known algorithm that combines resource packing and task dependencies. It yields an $O(\log n)$ -approx ratio on a DAG with n tasks [12]. The key idea is to partition tasks into *levels* such that all dependencies go across levels. The primary drawback with StripPart is that it prevents overlapping independent tasks that happen to be in different levels. A secondary drawback is that the recommended packers (e.g., [13]) do not support multiple resources and vector packing. We see that in practice StripPart under-performs the simpler heuristics.

8.4 How close is DAGPS to Optimal?

Figure 13 compares DAGPS with NewLB and the best previous lower bound $\max(\text{CPLen}, \text{TWork})$ (see § 4.2). Since the optimal schedule is no shorter than the lower bound, the figure shows that DAGPS is optimal for about 74% of DAGs. For half (three quarters) of the DAGs, DAGPS is within 32% of the new lower bound. A gap still remains: for the worst 13% of DAGs, DAGPS takes 47% longer. Manually examining these DAGs shows that NewLB is loose for most of them. However, the figure also shows that the NewLB improves upon previous lower bounds by almost 47% for most of the DAGs. We conclude that while more work remains towards a good lower bound, NewLB succeeds to argue that DAGPS is close to optimal for most of the production DAGs.

8.5 Sensitivity Analysis

We evaluate DAGPS’s sensitivity to parameter choices.

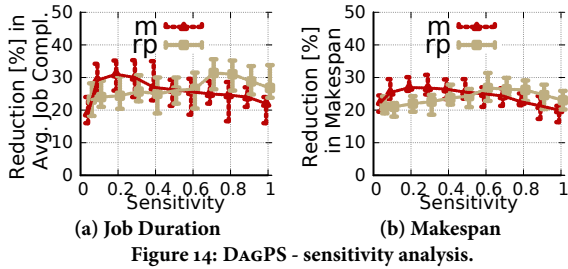


Figure 14: DAGPS - sensitivity analysis.

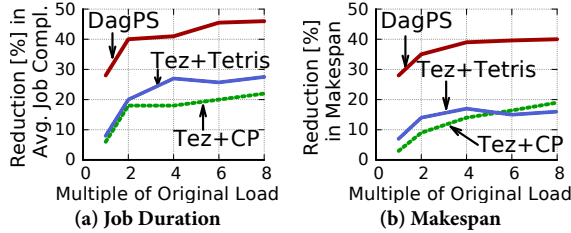


Figure 15: DAGPS's gains increase with cluster load.

Packing vs. Shortest Remaining Processing Time (srpt):

Recall that we combine packing score and `srpt` using a weighted sum with η (first box in Figure). Let η be m times the average over the two expressions that it combines. Here, we evaluate the sensitivity of the choice of m . Figure shows the reduction in average job completion time (on left) and makespan (on right) for different values of m . Values of $m \in [0.2, 0.4]$ have the most gains. Lower values lead to worse average job completion time because the effect of `srpt` reduces. On the other hand, larger values lead to moderately worse makespan. Hence, we recommend $m = 0.2$.

Remote Penalty: DAGPS uses a remote penalty `rp` to prefer local placement. Our analysis shows that both job completion time and makespan improve the most when `rp` is between 0.2 and 0.4 (Fig.). Since `rp` is a multiplicative penalty, lower values of `rp` cause the scheduler to miss (non-local) scheduling opportunities whereas higher `rp` can over-use remote resources on the origin servers. We use `rp = 0.2`.

Cluster Load: We vary cluster load by reducing the number of available servers without changing the workload. Figure shows the job completion times and makespan for a query set derived from TPC-DS. We see that both DAGPS and the alternatives offer more gains at higher loads. This is to be expected since the need for careful scheduling and packing increases when resources are scarce. Gains due to DAGPS increase by +10% at $\times 2$ load and by +15% at $\times 4$ load. However, the gap between DAGPS and the alternatives remains similar across load levels.

9. APPLYING DAGPS TO OTHER DOMAINS

We evaluate DAGPS's effectiveness in scheduling the DAGs that arise in distributed compilation systems [1, 2] and request-response workflows for Internet services [3].

Distributed build systems speed up the compilation of large code bases [4, 5]. Each build is a DAG with dependencies between the various tasks (compilation, linking, test, code analysis). The tasks have different runtimes and have

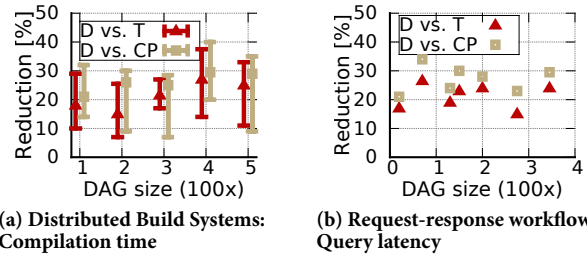


Figure 16: Comparing DAGPS (D) with Tetris (T) and Critical path scheduling (CP) on DAGs from two other domains.

different resource profiles. Figure a shows that DAGPS is faster than Tetris and faster than CP when scheduling the build DAGs from a production distributed build system. Each bar shows the median gain for DAGs of a certain size and the error bars are quartiles. The gains hold across DAG sizes/types.

We also examine the DAGs that arise in datacenter-side request-response workflows for Internet-services [6]. For instance, a search query translates into a workflow of dependent RPCs at the datacenter (e.g., spell check before index lookup, video and image lookup in parallel). The RPCs use different resources, have different runtimes and often execute on the same server pool [7]. Over several workflows from a production service, Figure b shows that DAGPS improves upon alternatives by about 20%.

These early results, though preliminary, are encouraging and demonstrate the generality of our work.

10. RELATED WORK

To structure the discussion, we ask four questions: (Q1) does a scheme consider both packing and dependencies, (Q2) does it make realistic assumptions, (Q3) is it practical to implement in cluster schedulers and, (Q4) does it consider multiple objectives such as fairness? To the best of our knowledge, DAGPS is unique in positively answering these four questions.

Q1: NO. Substantial prior work ignores dependencies but packs tasks with varying demands for multiple resources [8, 9, 10, 11, 12]. The best results are when the demand vectors are *small* [13]. Other work considers dependencies but assumes homogeneous demands [14, 15]. A recent multi-resource packing scheme, Tetris [16], succeeds on the three other questions but does not handle dependencies. Hence, we saw in §3 that it performs poorly when scheduling DAGs. Further, Tetris has poor worst-case performance (up to 100 times slower, see Figure 3) and can be arbitrarily unfair.

Q1: YES, Q2: NO. The packing+dependencies problem has been considered at length under the keyword *job-shop scheduling* [17, 18, 19, 20]. Most results assume that jobs are known a priori (i.e., the offline case). See [21] for a survey. For the online case (the version considered here), no algorithms with bounded competitive ratios are known [22, 23]. Some other notable work assumes only two resources [24], applies only for a chain but not a general DAG [25] or assumes one

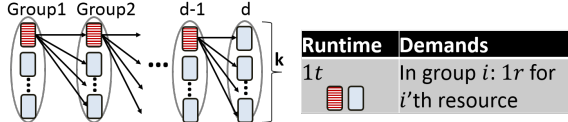


Figure 17: A counter-example DAG that shows any scheduler not considering DAG structure will be (d) times OPT.

cluster-wide resource pool [10].

Q3 : NO. All of the schemes listed above consider one DAG at a time and are not easily adaptable to the online case when multiple DAGs share a cluster. Work on related problems such as VM allocation [11] also considers multi-resource packing. However, cluster schedulers have to support roughly two to three orders of magnitude higher rate of allocation (tasks are more numerous than VMs).

Q3 : YES, Q1 : NO. Several notable works in cluster scheduling exist such as Quincy [12], Omega [13], Borg [14], Kubernetes [15] and Autopilot [16]. None of these combine multi-resource packing with DAG-awareness. Many do neither. Job managers such as Tez [17] and Dryad [18] use simple heuristics such as breadth-first scheduling which perform quite poorly in our experiments.

Q4 : NO. There has been much recent work on novel fairness schemes to incorporate multiple resources [19] and be work-conserving [20]. Several applications arise especially in scheduling co-flows [21, 22]. We note that these fairness schemes neither pack nor are DAG-aware. DAPS can incorporate these fairness methods as one of the multiple objectives and trade off bounded unfairness for performance.

11. CONCLUDING REMARKS

DAGs are indeed a common scheduling abstraction. However, we found that existing algorithms make several key assumptions that do not hold in practical settings. Our solution, DAPS is an efficient online solution that scales to large clusters. We experimentally validated that it substantially improves the scheduling of DAGs in both synthetic and emulated production traces. The core contributions are three-fold: (1) constructing a good schedule by placing tasks out-of-order on to a virtual resource-time space, (2) an online heuristic that solely enforces the desired schedules and helps with other concerns such as packing and fairness, and (3) an improved lower bound that lets us show that our heuristics are close to optimal. Much of these innovations use the fact that job DAGs consist of groups of tasks (in each stage) that have similar durations, resource needs and dependencies. We intend to contribute our DAPS implementation to Apache YARN/Tez projects. As future work, we are considering applying these DAG scheduling ideas to related domains, most notably scheduling the co-flows with dependencies that arise in geo-distributed analytics [23, 24, 25].

A. VALUE OF DAG AWARENESS

Proof of Lemma 1: Figure 17 shows an adversarial DAG for which any scheduler that ignores dependencies will take

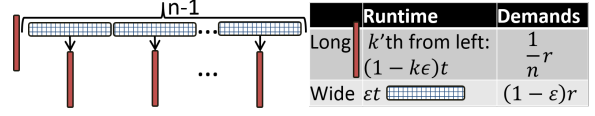


Figure 18: An example DAG where critical path scheduling is $O(n)$ times OPT where n is the number of nodes in the DAG.

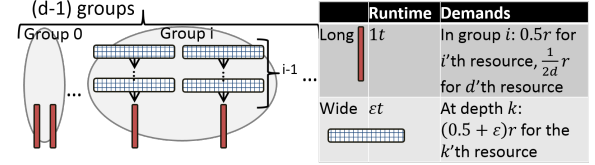


Figure 19: An example DAG where Tetris [37] is $d -$ times OPT when tasks use d kinds of resources.

(d) times OPT. Assume cluster capacity is r for each of the d resources. The DAG has d groups, each having a task (indicated with red dashes) that is the parent of all the tasks in the next group. This information is unavailable (and unused) by schedulers that do not consider the DAG structure. Hence, regardless of which order the scheduler picks tasks, an adversary can choose the last task in a group to be the red one. Hence, such schedulers will take kdt time. OPT only requires $(k + d - 1)t$ since it can schedule the red tasks first (in $(d - 1)t$) and afterwards one task from each group can run simultaneously (kt more steps). We use Yao's max-min principle [26] (the lower bound on any randomized algorithm is the same as lower bound on deterministic algorithms with randomized input) to extend the counter-example. If we randomized the choice of the red task, the expected time at which the red task will finish is kt/d and hence the expected schedule time is $k(d + 1)t/d$ which is still (d) times OPT. \square

B. WORST-CASE DAG EXAMPLES

Proof of Lemma 2: Figure 18 shows an example DAG where CPSched takes n times worse than OPT for a DAG with n tasks. The critical path lengths of the various tasks are such that CPSched alternates between one long task and one wide task left to right. However, it is possible to overlap all of the long running tasks. This DAG completes at $\sim nt$ and $\sim t$ with CPSched and OPT respectively. Fig. 19 shows an example where Tetris [37] is $d -$ times OPT. As in the above example, all long tasks can run together, hence OPT finishes in t . Tetris greedily schedules the task with the highest dot-product between task demands and available resources.

The DAG is constructed such that whenever a long task is runnable, it will have a higher score than any wide task. Further, for every long task that is not yet scheduled, there exists at least one wide parent that cannot overlap any long task that may be scheduled earlier. Hence, Tetris takes $(d - 1)t$ which is $(d - 1)$ times OPT. Combining these two principles, we conjecture that it is possible to find similar examples for any scheduler that ignores dependencies or is not resource-aware. \square

12. REFERENCES

- [] bigdata platforms and bigdata analytics software. <http://bit.ly/1DR0qgt>.
- [] Apache Hadoop. <http://hadoop.apache.org>.
- [] Apache Tez. <http://tez.apache.org/>.
- [] Bazel. <http://bazel.io/>.
- [] Big-Data-Benchmark. <http://bit.ly/1H1FRH0>.
- [] Condor. <http://research.cs.wisc.edu/htcondor/>.
- [] Hadoop: Fair scheduler/ slot fairness. <http://bit.ly/1PfsT7F>.
- [] Hadoop YARN Project. <http://bit.ly/1iS8xvP>.
- [] Kubernetes. <http://kubernetes.io/>.
- [] Market research on big-data-as-service offerings. <http://bit.ly/1V6TXV4>.
- [] TPC-DS Benchmark. <http://bit.ly/1J6uDap>.
- [] TPC-H Benchmark. <http://bit.ly/1KRK5gl>.
- [] Rate control for communication networks: shadow prices, proportional fairness and stability. *Journal of the Operational Research Society* ().
- [] A. G , N. J , S. K , C. K , P. L , D. A. M , P. P , S. S . VL : A Scalable and Flexible Data Center Network. In *SIGCOMM* ().
- [] A , S., K , S., B , N., W , M.-C., S , I., Z , J. Re-optimizing data parallel computing. In *NSDI* ().
- [] A -F , M., L , A., V , A. A scalable, commodity data center network architecture. In *SIGCOMM* ().
- [] A , G., . Reining in the Outliers in MapReduce Clusters Using Mantri. In *OSDI* ().
- [] A , E., B , D., C , K., K , T., S , N., S , C., S , R., T , R., W , J., Z , Y. Value-maximizing deadline scheduling and its application to animation rendering. In *SPAA* ().
- [] A , M., . Spark sql: Relational data processing in spark. In *SIGMOD* ().
- [] A , J., B , S., I , S. Strip packing with precedence constraints and strip packing with release times. In *SPAA* (), ACM.
- [] A , Y., C , I. R., F , A., R , A. Packing small vectors. In *SODA* ().
- [] B , K. P., B , P. An approximate algorithm for the partitionable independent task scheduling problem. *Urbana 51* (), .
- [] B , N., J , S., Z , J. Continuous cloud-scale query optimization and processing. *VLDB* ().
- [] C , R., . SCOPE: Easy and Efficient Parallel Processing of Massive Datasets. In *VLDB* ().
- [] C , C., . Flumejava: easy, efficient data-parallel pipelines. In *PLDI* ().
- [] C , C., S , K. On multidimensional packing problems. *SIAM J. Comput.* ().
- [] C , M., L , Z., G , A., S , I. Hug: Multi-resource fairness for correlated and elastic demands. In *NSDI* ().
- [] C , M., Z , Y., S , I. Efficient co ow scheduling with varies. In *SIGCOMM* ().
- [] C , N., R , M., B , R. Virtual Network Embedding with Coordinated Node and Link Mapping. In *INFOCOM* ().
- [] C , E.G., J., G , R. Optimal scheduling for two-processor systems. *Acta Informatica* ().
- [] C , A., S , C. A new algorithm approach to the general Lovász local lemma with applications to scheduling and satisfiability problems (extended abstract). In *STOC* ().
- [] D , J., G , S. Mapreduce: Simplified data processing on large clusters. In *OSDI* ().
- [] G , A., . Dominant Resource Fairness: Fair Allocation Of Multiple Resource Types. In *NSDI* ().
- [] G , M., S , W., P , C., V , D., N , I., L , B. Automated migration of build scripts using dynamic analysis and search-based refactoring. In *OOPSLA* ().
- [] G , L. A., P , M., S , A., S , E. Better approximation guarantees for job-shop scheduling. In *SODA* ().
- [] G , R. L. Bounds on multiprocessing timing anomalies. *SIAM Journal on Applied Mathematics* ().
- [] G , R., A , G., K , S., R , S., A , A. Multi-resource Packing for Cluster Schedulers. In *SIGCOMM* ().
- [] H , B., K , A., Z , M., G , A., J , A. D., K , R., S , S., S , I. Mesos: a platform for fine-grained resource sharing in the data center. In *Proceedings of the 8th USENIX conference on Networked systems design and implementation* (Berkeley, CA, USA,), NSDI' , USENIX Association, pp. - .
- [] H , C.-C., G , L., Y , M. Scheduling jobs across geo-distributed datacenters. In *SOCC* ().
- [] H , C.-C., G , L., Y , M. Scheduling jobs across geo-distributed datacenters. In *SOCC* ().
- [] I , M. Autopilot: Automatic Data Center Management. *OSR* ().
- [] I , M., . Dryad: Distributed Data-parallel Programs from Sequential Building Blocks. In *Eurosys* ().
- [] I , M., . Quincy: Fair Scheduling For Distributed Computing Clusters. In *SOSP* ().

- [] J. R. C. D., H. W. A quantitative measure of fairness and discrimination for resource allocation in shared computer systems. *CoRR cs.NI/9809099* ().
- [] J. V., B. P., K. S., M. I., R. M., Y. C. Speeding up distributed request-response work flows. In *SIGCOMM* ().
- [] K. Q., I. M., Y. Y. Optimus: A dynamic rewriting framework for data-parallel execution plans. In *Eurosys* ().
- [] K. Y., A. I. Static scheduling algorithms for allocating directed task graphs to multiprocessors. *ACM Computing Surveys (CSUR)* ().
- [] L. F. T., M. B. M., R. S. Universal packet routing algorithms. In *FOCS* ().
- [] L. R., T. D., W. G. J. Approximation Algorithms for Scheduling Malleable Tasks Under Precedence Constraints. *International Journal of Foundations of Computer Science* ().
- [] M. Z., A. K., A. D. J., R. K., I. S. Improving MapReduce Performance in Heterogeneous Environments. In *OSDI* ().
- [] M. B. M., H. F. M., V. B., W. M. Exploiting locality for data management in systems of limited bandwidth. In *FOCS* ().
- [] M., M., S., O. (acyclic) job shops are hard to approximate. In *FOCS* ().
- [] M., M., O., S. Improved bounds for flow shop scheduling. In *ICALP* ().
- [] M. C., I. S. Efficient flow scheduling without prior knowledge. In *SIGCOMM* ().
- [] M., R., R., P. *Randomized Algorithms*. Cambridge University Press, New York, NY, USA, .
- [] O., K., W., P., Z., M., S., I. Sparrow: Distributed, low latency scheduling. In *SOSP* ().
- [] P., R., . Heuristics for Vector Bin Packing. In *MSR TR* ().
- [] P., Q., A., G., B., P., K., S., A., A., B., P., S., I. Low latency geo-distributed analytics. In *SIGCOMM* ().
- [] R., A., Z., H., B., J., P., G., S., A. C. Inside the social network's (datacenter) network. In *SIGCOMM* ().
- [] S., I. Reverse-fit: A ϵ -optimal algorithm for packing rectangles. In *Proceedings of the Second Annual European Symposium on Algorithms* ().
- [] S., E., B., J. The User and Business Impact of Server Delays, Additional Bytes, and Http Chunking in Web Search. <http://velocityconf.com/velocity2009/public/schedule/detail/8523>, .
- [] S., M., K., A., A. E-M., M., W., J. Omega: Flexible, scalable schedulers for large compute clusters. In *EuroSys* ().
- [] S., D. B., S., C., W., J. Improved approximation algorithms for shop scheduling problems. *SIAM J. Comput.* ().
- [] S., M., . Efficient fair queueing using deficit round robin. In *SIGCOMM* ().
- [] S., A., . Jupiter rising: A decade of clos topologies and centralized control in google's datacenter network. In *SIGCOMM* ().
- [] S., I., N., K., J., K., D., P. Tight bounds for online vector scheduling. In *FOCS* ().
- [] S., L., C., M., S., S., F., A. C.: Cutting tail latency in cloud data stores via adaptive replica selection. In *NSDI* ().
- [] T., A., . Hive- a warehousing solution over a map-reduce framework. In *VLDB* ().
- [] V., A., P., L., K., M. R., O., D., T., E., W., J. Large-scale cluster management at Google with Borg. In *EuroSys* ().
- [] V., A., C., C., G., P. B., J., T., P., J., V., G. Global analytics in the face of bandwidth and regulatory constraints. In *NSDI* ().
- [] W., G. J. There Is No Asymptotic PTAS For Two-Dimensional Vector Packing. In *Information Processing Letters* ().
- [] Y., A., I., C., S., K., B., S. Tight bounds for online vector bin packing. In *STOC* ().
- [] Y., A., I., R., C., I., G. The loss of serving in the dark. In *STOC* ().
- [] Z., M., . Spark: Cluster computing with working sets. Tech. Rep. UCB/EECS- , EECS Department, University of California, Berkeley, .