



HAL
open science

Data-Driven Locality-Aware Batch Scheduling

Maxime Gonthier, Elisabeth Larsson, Loris Marchal, Carl Nettelblad, Samuel Thibault

► **To cite this version:**

Maxime Gonthier, Elisabeth Larsson, Loris Marchal, Carl Nettelblad, Samuel Thibault. Data-Driven Locality-Aware Batch Scheduling. APDCM 2024 - 26th Workshop on Advances in Parallel and Distributed Computational Models, 38th IEEE International Parallel and Distributed Processing Symposium, May 2024, San Francisco, United States. hal-04500281

HAL Id: hal-04500281

<https://inria.hal.science/hal-04500281>

Submitted on 11 Mar 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Data-Driven Locality-Aware Batch Scheduling

Maxime Gonthier
University of Chicago, ENS de Lyon & INRIA
Department of computer science
Chicago, USA
mgonthier@uchicago.edu

Elisabeth Larsson
Uppsala University
Department of Information Technology
Uppsala, Sweden
elisabeth.larsson@it.uu.se

Loris Marchal
CNRS
ENS de Lyon
Lyon, France
loris.marchal@ens-lyon.fr

Carl Nettelblad
Uppsala University
Department of Information Technology
Uppsala, Sweden
carl.nettelblad@it.uu.se

Samuel Thibault
University of Bordeaux
LaBRI
Talence, France
samuel.thibault@inria.fr

Abstract—Clusters employ workload schedulers such as the Slurm Workload Manager to allocate computing jobs onto nodes. These schedulers usually aim at a good trade-off between increasing resource utilization and user satisfaction (decreasing job waiting time). However, these schedulers are typically unaware of jobs sharing large input files, which may happen in data intensive scenarios. The same input files may end up being loaded several times, leading to a waste of resources.

We study how to design a *data-aware job scheduler* that is able to keep large input files on the computing nodes, without impacting other memory needs, and can benefit from previously-loaded files to *decrease data transfers in order to reduce the waiting times of jobs*.

We present three schedulers capable of distributing the load between the computing nodes as well as re-using input files already loaded in the memory of some node as much as possible.

We perform simulations with single node jobs using traces of real HPC-cluster usage, to compare them to classical job schedulers. The results show that keeping data in local memory between successive jobs and using data-locality information to schedule jobs improves performance compared to a widely-used scheduler (FCFS, with and without backfilling): a reduction in job waiting time (a 7.5% improvement in stretch), and a decrease in the amount of data transfers (7%).

Index Terms—Batch scheduling, Job input sharing, Data-aware, Job scheduling, High Performance Data Analytics

I. INTRODUCTION

To meet the ever-increasing demand for scientific computation power, High-Performance Computing platforms are typically composed of large numbers of computation nodes. Scientists typically submit their computation jobs to a scheduler, which decides the ordering and mapping of the jobs on the platform. This needs to be performed with particular care of balancing between resource utilization and user satisfaction, so as to leverage the computation resources as efficiently as possible, while avoiding adverse pathological cases that could particularly impact some users rather than others.

Computation jobs however need data input which, more often than not, can be very large, notably for many sub-fields of life science with highly data-dependent workflows like taxonomic identification of DNA fragments or ancestral

reconstructions. It is a frequent use pattern for users of such communities to submit large batches of jobs using the same input files. Loading such data input from the storage nodes may consume a significant part of the job duration. This load penalty can however be avoided altogether when the data was actually already used by the previous job running on the computation node, and thus still immediately available on the node. Taking care of scheduling jobs that use the same data input one after the other thus allows to reduce the jobs completion times, leading to better platform usage efficiency. Unfortunately, classical job schedulers mostly do not take data input into account, and thus do not benefit from such data reuse; most jobs always have to re-load their data input.

In this paper, we propose to model the benefits of re-using data loads between successive jobs, and we introduce new algorithms that add such data reuse to the scheduling balance. By tracking which data is loaded on which node for the scheduled jobs, they are able to significantly reduce data loads, thus improving both resource utilization and user satisfaction. We evaluated these algorithms thanks to traces of actual jobs submissions observed on a large cluster platform. This allows to assess the effectiveness of our schedulers over a variety of realistic working sets. This revealed that while our schedulers get slightly worse results over some working set samples (those which exhibit cluster ample underuse), most working set samples largely benefit from our schedulers, leading to interesting benefit overall.

We thus present the following contributions in this paper:

- 1) We formalize our model of scheduling data-intensive jobs sharing input files on a cluster (Section III).
- 2) We propose three new schedulers focusing on re-using input files while minimizing evictions and avoiding starvation (Section IV).
- 3) We extract job information from historical logs of an HPC-cluster to build workloads that correspond to the needs and behaviors of real users (Section V-B).
- 4) We implement all three schedulers as well as two state-of-the-art schedulers on a simulator and study the perfor-

mance (mean stretch, amount of time spent waiting for a file to be loaded and number of core hours) obtained after scheduling almost 2 millions jobs (Section V). Our evaluation demonstrates that our schedulers in most cases surpass the state-of-the-art schedulers. We show that workloads that heavily saturate the cluster benefit much more from our strategies, which results in considerable reduction in the job waiting times.

We used simulations to avoid disrupting users on a production cluster, but our strategies could be implemented on the Slurm Workload Manager by asking the users to flag their input files.

II. RELATED WORK

A. Scheduling jobs on large clusters

A various number of workloads managers have emerged as a way to meet the rising numbers of high-performance computing clusters. Workload managers like Slurm [1], OAR [2], TORQUE [3], LoadLeveler [4], Portable Batch System [5] or SunGrid [6] all offer various scheduling strategies.

The First-Come-First-Served (FCFS) algorithm is the prevalent default scheduler on most of these solutions [7]. Moreover, Slurm is used on most of the TOP500 supercomputers and its default strategy is FCFS [8] as well. Backfilling strategies are known to improve the use of supercomputer resources [9]. The most commonly-used backfilling strategy, called conservative backfilling [10] [11], follows a simple paradigm: "a job can only be backfilled if it does not delay a job of higher priority". We can then safely assume that comparing ourselves to FCFS with and without conservative backfilling can bring significant insights on what improvements can be achieved on data-intensive workloads.

Other scheduling strategies exist like Maui [12], Gang scheduling [13], RASA [14], that use the advantages of both Min-min and Max-min algorithms, RSDC [15] that divides large jobs in sub-jobs for a more refined scheduling, or PJSC [16] that is a priority-based scheduler; however these heuristics do not consider the impact that input re-use could have on data-intensive workloads. We aim at resolving this issue in this paper.

B. Using distributed file systems to deal with data-intensive workloads

Distributed file systems are a solution to ease the access to shared input files. They facilitate the execution of I/O-intensive batch jobs by selecting appropriate storage policies. BDFS [17] is designed to orchestrate large, I/O-intensive batch workloads on remote computing clusters. HDFS [18] (Hadoop Distributed File System) incorporates storage-aware scheduling. It migrates a computation closer to where the data is located, rather than moving the data to where the application is running, in order to reduce communication. These solutions are mainly storage systems that use a history of file locations to serve as a backup. In our scenario, we copy the data from an already-redundant system (an online database for example) and store it locally on the node in an ephemeral way. Thus, in the event of a crash, we do not manage the data which is already

redundant, it simply results in an aborted job. Secondly, the scheduling can cause issues (notably MapReduce, used in HDFS), as described by Weets et al. [19]. By not using a distributed file system, we avoid these problems altogether. Lastly, file systems are particularly efficient when the input data used are identical over time. In our case, between users, the inputs will be largely different, making distributed file systems less efficient.

C. Using schedulers to deal with data-intensive workloads

Some schedulers tackle the issue of data-intensive workloads. A solution can be to minimize network contention by allocating nodes to even out node and switch contention [20]. In our model, we are not studying the network topology and consider independent nodes. This is reasonable, since our main concern is the cross-section bandwidth to a shared storage solution. Nikolopoulos et al. [21] focus on a better utilization of idling memory together with thrashing limitation. Our focus will be to control data loads in order to limit eviction and will thus naturally limit thrashing. Agrawal et al. [22] propose to schedule jobs not sharing a file first and to use a stochastic model of job arrivals for each input file to maximize re-use. That work is aimed at the Map-Reduce model and allows to predict future jobs arrivals, two prerequisites that we do not consider.

III. FRAMEWORK

We consider the problem of scheduling a set of \mathcal{J} independent jobs, denoted $\mathbb{J} = \{J_1, J_2, \dots, J_{\mathcal{J}}\}$ on a set of \mathcal{N} nodes: $\mathbb{N} = \{\text{Node}_1, \text{Node}_2, \dots, \text{Node}_{\mathcal{N}}\}$. Each node Node_i is equipped with m cores noted: c_1^i, \dots, c_m^i sharing a memory of size M .

Each job J_i depends on an input file noted $File(J_i)$, which is initially stored in the main shared file system. During the processing of a job J_i on Node_k , $File(J_i)$ must be in the memory of Node_k . If this is not the case before starting computation of job J_i , then $File(J_i)$ is loaded into the memory. We denote by $\mathbb{F} = \{F_1, F_2, \dots, F_{\mathcal{F}}\}$ the set of distinct input files, whose size is denoted by $Size(F_i)$. Each job runs on a single node, but they can make use of a different number of cores. Each job J_i has the following attributes:

- Resource requirement: job J_i requests $Cores(J_i)$ cores, such that $1 \leq Cores(J_i) \leq m$;
- Input file: $File(J_i) \in \mathbb{F}$;
- Submission date: $SubDate(J_i)$;
- Requested running time: $WallTime(J_i)$: if not finished after this duration, job J_i is killed by the scheduler;
- Actual running time: $Duration(J_i)$ (unknown to the scheduler before the job completion).

Each of the \mathcal{J} jobs must be processed on one of the \mathcal{N} nodes. As stated earlier, the shared file system initially contains all files in \mathbb{F} . Each node is connected to the file system with a link of limited bandwidth, denoted by $Bandwidth$: transferring a data of size S from the shared file system to some node's memory takes a time $S/Bandwidth$. The limited bandwidth

as well as the large file sizes are the reasons why we aim at restricting the amount of data movement.

Data-intensive jobs generally have very scarce outputs compared to the large inputs they use (e.g., genome databases). Therefore, we do not consider the output of jobs. We assume that jobs are devoted a fraction of the memory of a node proportional to the number of requested cores, so that jobs willing to process large input files must request large number of cores. This way, we make sure that the memory of a node is large enough to accommodate all input files of running jobs. A file stored in the memory of a node can be shared by two jobs J_i and J_j only in either of the following situations:

- 1) J_i and J_j are computed in parallel on the same node.
- 2) J_i and J_j are computed on the same node consecutively (i.e., no job is started on this node between the completion of J_i and the start of J_j).

This can hold true if the file data is accessed through I/O (traditional or memory-mapped), allowing the same page cache to serve multiple processes from different jobs. Otherwise we consider that memory operations of jobs scheduled between J_i and J_j will cause the file to be evicted.

For each job J_i , the scheduler is in charge of deciding which node will process J_i , and which cores of this node are allocated to the job, as well as a starting time t_k for J_i . Job J_i is thus allocated a time window from t_k to $t_k + WallTime(J_i)$ devoted to (i) possibly loading the input file $File(J_i)$ in the memory (if it is not already present at time t_k) and (ii) executing job $J(i)$. If the job is not completed at time $t_k + WallTime(J_i)$, it is killed by the scheduler to ensure that later jobs are not delayed. The scheduler must also make sure that no two jobs are executed simultaneously on the same cores.

It is important to note that the file transfer is done before the computation and cannot be overlapped. Jobs are non-preemptible: when started, a job is executed throughout its completion.

Our objective is to reach an efficient usage of the platform and to limit job waiting times. Each user submitting jobs is interested in obtaining the result of jobs as soon as possible. Hence we focus on the time spent in the system for each job, also called the flow time (or flow) or turnaround time of the job:

$$Flow(J_i) = CompletionDate(J_i) - SubDate(J_i).$$

In the following, we want to consider aggregated performance metrics on job flows, such as average flow. However, the duration of a job significantly impacts its flow time. Jobs with the same flow but very different durations do not experience the same quality of service. To avoid this, the *stretch* metric has been introduced that compares the actual flow of a job to the one it would experience on an empty cluster:

$$ReferenceFlow(J_i) = \frac{Size(File(J_i))}{Bandwidth} + Duration(J_i)$$

$$stretch(J_i) = \frac{Flow(J_i)}{ReferenceFlow(J_i)}.$$

The stretch represents the slow-down of a job due to sharing the platform with other jobs and is a commonly used metric [23]. Considering the stretch allows to better aggregate performance from small and large jobs.

IV. JOB SCHEDULING WITH INPUT FILES

Here, we present various schedulers used to allocate jobs to computing resources. We start with two reference schedulers (FCFS and EFT) and then move to proposing three locality-aware job schedulers (named LEA, LEO and LEM). Each of these five schedulers can be used both without or with backfilling. For the sake of clarity, we first present the simpler version, without backfilling, before detailing the modifications needed to include backfilling.

As presented above, the task of the scheduler is to allocate a set of cores of one node to each job submitted until now: some jobs may be started right away, while other jobs may be delayed and scheduled later: resource reservations are made for these jobs. Jobs are presented to the scheduler in the form of a global queue, sorted by job submission time. These scheduling policies are online algorithms which are called each time a task completes (making cores available) or upon the submission of a new job.

A. Two schedulers from the state of the art: FCFS and EFT

A widely-used job scheduler that is typically considered to be efficient is First-Come-First-Serve (FCFS), detailed in Algorithm 1. Implementing this scheduler requires to remember the time of next availability for each core. Then, for each job, we look for the first time when a sufficient number of cores is available, and we allocate the job to those cores.

Algorithm 1 First-Come-First-Serve (FCFS)

- 1: **for each** $J_i \in$ the jobs queue **do**
 - 2: **for each** $Node_k \in \mathbb{N}$ **do**
 - 3: Find smallest time t_k such that $Cores(J_i)$ cores are available on $Node_k$
 - 4: Select $Node_k$ with the smallest t_k
 - 5: Schedule J_i on $Cores(J_i)$ cores of $Node_k$ that are available starting from t_k
 - 6: Mark these cores busy until time $t_k + WallTime(J_i)$
-

FCFS is a standard baseline comparison for job scheduling. However, it is not aware of the capability of the system to keep a large data file in the memory of a node between the execution of two consecutive jobs. A first step towards a locality-aware scheduler is to select a node for each job not only based on the cores' availability time, but also using the file availability time, based on the file transfer time. This is the purpose of the Earliest-Finish-Time (or EFT) scheduler, described in Algorithm 2: by selecting the node that can effectively start the job at the earliest time, it minimizes the job completion time. There are three scenarios to compute the time t'_k at which the input file of a job J_i is available on $Node_k$:

- 1) A job started before t_k on $Node_k$ uses the same input file and it is already in memory, thus $t'_k = t_k$;

- 2) The input file of J_i is not in memory, thus $t'_k = t_k + \frac{Size(File(J_i))}{Bandwidth}$;
- 3) The input file of J_i is partially loaded on $Node_k$: this happens when some job J_j using the same input file has been scheduled on other cores of the same node at time $StartTime(J_j) < t_k$ but the file transfer has not been completed at time t_k . Then the file will be available at time: $t'_k = StartTime(J_j) + \frac{Size(File(J_i))}{Bandwidth}$.

Algorithm 2 Earliest-Finish-Time (EFT)

- 1: **for each** $J_i \in$ the jobs queue **do**
 - 2: **for each** $Node_k \in \mathbb{N}$ **do**
 - 3: Find smallest time t_k such that $Cores(J_i)$ cores are available on $Node_k$
 - 4: Find time $t'_k \geq t_k$ at which $File(J_i)$ is available on $Node_k$
 - 5: Select $Node_k$ with the smallest t'_k
 - 6: Schedule J_i on $Cores(J_i)$ cores of $Node_k$ that are available starting from t_k
 - 7: Mark these cores busy until time $t_k + WallTime(J_i)$
-

B. Data-locality-based schedulers

The previous strategies focus on starting (FCFS) or finishing (EFT) a job as soon as possible, respectively. Those are good methods to avoid node starvation and reduce queue times. However, they may lead to loading the same input file on a large number of nodes in the platform, only to minimize immediate queue times. Time is thus spent loading the input file multiple times. This can affect the global performance of the system by delaying subsequent jobs. We present three strategies that attempt to take data locality into account in a better way to reduce queuing times in the long run by increasing data reuse.

The first proposed algorithm, called Locality and Eviction Aware (LEA) and detailed in Algorithm 3, aims at a good balance between node availability and data locality. We consider three quantities to rank nodes:

- The availability time for computation t_k ;
- The time needed to complete loading the input file for the job on $Node_k$ ($t'_k - t_k$);
- The time required to reload files that need to be evicted before loading the input file; this time is computed using all files in memory and considering that a fraction of these files need to be evicted, corresponding to the fraction of the memory used by the job.

The intuition for the third criterion is that if loading a large file in memory requires the eviction of many other files, these files will not be available for later jobs and may have to be reloaded. In the LEA strategy, we put a strong emphasis on data loading, in order to really favor data locality. Hence, when computing the score for each node $Node_k$, we sum the previous three quantities, with a weight W for the second one (loading time). In our experiments, based on empirical evaluation, we set this value to $W = 500$, incidentally roughly

Algorithm 3 Locality and Eviction Aware (LEA)

- 1: **for each** $J_i \in$ the jobs queue **do**
 - 2: **for each** $Node_k \in \mathbb{N}$ **do**
 - 3: Find smallest time t_k such that $Cores(J_i)$ cores are available
 - 4: Find time $t'_k \geq t_k$ at which $File(J_i)$ is available on $Node_k$
 - 5: $LoadOverhead \leftarrow t'_k - t_k$
 - 6: Let \mathfrak{F} be the set of files in the memory of node $Node_k$ at time t_k
 - 7: $EvictionPenalty \leftarrow (\sum_{F_j \in \mathfrak{F}} Size(F_j) \times Size(File(J_i)) / M) / Bandwidth$
 - 8: $score_{Node_k} \leftarrow t_k + W \times LoadOverhead + EvictionPenalty$
 - 9: Select $Node_k$ with the smallest $score_{Node_k}$
 - 10: Schedule J_i on $Cores(J_i)$ cores of $Node_k$ that are available starting from t_k
 - 11: Mark these cores busy until time $t_k + WallTime(J_i)$
-

equivalent to the number of nodes. Note that the other two quantities usually have very different values: the availability time is usually much larger than the time for reloading evicted data. Hence this last criterion is mostly used as a tie-break in case of equality of the first two criteria.

The LEA strategy puts a dramatic importance on data loads. Hence, it is very useful when the platform is fully loaded and some jobs can safely be delayed to favor data reuse and avoid unnecessary loads. However, when the platform is not fully loaded, delaying jobs can be detrimental, as it can increase the response time for some jobs, without any benefit for other jobs. Our second proposed strategy, named Locality and Eviction Opportunistic (LEO) and described in Algorithm 4, tries to adapt based on the current cluster load: if we find some nodes that can process the job right away, we select the one that will minimize the completion time (as in the EFT strategy). Otherwise, we assume that the platform is fully loaded and we apply the previous LEA strategy, to favor data reuse.

We present a third strategy, described in Algorithm 5, called Locality and Eviction Mixed (LEM). LEM takes a similar approach to LEO but performs a simple mix between the EFT and the LEA strategies: when the platform is saturated (each node is running at least one job), the LEA strategy is applied, otherwise the EFT strategy is used. The two main differences with LEO are that LEM's consideration of a saturated cluster is broader than LEO's, and that the strategy shift is applied to all jobs until the cluster is no longer saturated. LEO considers saturation on a per-job basis, allowing jobs that use a small number of cores to be placed on nodes occupied by a job that is already running but does not use the entire node. LEM considers a node to be saturated even if the job running on it is not using it completely. This is particularly useful when the job queue consists of jobs that require a full node.

Algorithm 4 Locality and Eviction Opportunistic (LEO)

```
1: for each  $J_i \in$  the jobs queue do
2:   for each  $\text{Node}_k \in \mathbb{N}$  do
3:     Find smallest time  $t_k$  such that  $\text{Cores}(J_i)$  cores
     are available
4:     Find time  $t'_k \geq t_k$  at which  $\text{File}(J_i)$  is available
     on  $\text{Node}_k$ 
5:     if  $t_k = \text{current\_time}$  then
6:        $\text{score}_{\text{Node}_k} \leftarrow t'_k$ 
7:     else
8:        $\text{LoadOverhead} \leftarrow t'_k - t_k$ 
9:       Let  $\mathfrak{F}$  be the set of files in the memory of node
      $\text{Node}_k$  at time  $t_k$ 
10:       $\text{EvictionPenalty} \leftarrow (\sum_{F_j \in \mathfrak{F}} \text{Size}(F_j) \times$ 
      $\text{Size}(\text{File}(J_i))/M)/\text{Bandwidth}$ 
11:       $\text{score}_{\text{Node}_k} \leftarrow t_k + W \times \text{LoadOverhead} +$ 
      $\text{EvictionPenalty}$ 
12:     Select  $\text{Node}_k$  with the smallest  $\text{score}_{\text{Node}_k}$ 
13:     Schedule  $J_i$  on  $\text{Cores}(J_i)$  cores of  $\text{Node}_k$  that are
     available starting from  $t_k$ 
14:     Mark these cores busy until time  $t_k + \text{WallTime}(J_i)$ 
```

Algorithm 5 Locality and Eviction Mixed (LEM)

```
1: for each  $J_i \in$  the jobs queue do
2:   if each node is running at least one job at
      $\text{current\_time}$  then
3:      $\text{LEA}(J_i, \mathbb{N})$ 
4:   else
5:      $\text{EFT}(J_i, \mathbb{N})$ 
```

C. Adding backfilling to all strategies

As mentioned above, backfilling has been proposed to increase the performance of production cluster job schedulers, by allowing jobs with lower priority to be scheduled before jobs with higher priority. In our setting, the priority is directly linked to the submission order: if J_i is submitted before J_j , then J_i has a higher priority than J_j . In order to avoid jobs being perpetually delayed, bounds have to be set on how already-scheduled jobs can be affected by backfilling. As discussed above, conservative backfilling is the most restrictive version and one of the most commonly-used strategies to improve cluster utilization. It forbids any modification on the resource reservations of high-priority jobs: a job may be scheduled before other jobs that appear earlier in the queue, provided that it does not impact the starting time of these jobs.

For each of the previous scheduling strategies, we consider a variant using conservative backfilling (suffixed by -BF). To add backfilling, Algorithms 1, 2, 3 and 4 have to be modified: we change the choice of the earliest time when resources are available for a job (Line 3). Instead of considering the time at which cores are (indefinitely) available, we look for an availability time window starting at t_k that is long enough to hold the job. Specifically, we change Line 3 into:

3': Find smallest time t_k such that $\text{Cores}(J_i)$ cores are available from t_k until $t_k + \text{WallTime}(J_i)$ on Node_k

V. PERFORMANCE EVALUATION AND ANALYSIS

Actual workloads on HPC resources shared by a great number of users with diverse needs can contain projects that might launch a burst of jobs using the same file just a few thousand core hours in length, then be quiet for a long time processing the results, and then launch another such burst. On our cluster, we have observed such behavior. Changing the scheduling strategy of a production cluster for scheduling research would disrupt the community using this cluster. For this reason, we choose to perform simulations based on logs of a real computing platform. In this section, we describe how we used these logs as well as the results of the corresponding simulations.

A. Platform description

We use Rackham, the main cluster for HPC workloads within the former Swedish National Infrastructure for Computing. It served half of all PIs there in a single year, and at least three Nature papers in 2023 acknowledged the use of Rackham. A variety of applications including genomics analysis are performed on this cluster. At the time of our analysis, it contained 486 nodes connected by Infiniband to a 6 PB Lustre storage system. Each node has two 10-core Intel Xeon V4 CPU at 2.20 GHz/core. Most of the nodes have 128 GB of RAM, and a few are larger. To avoid an additional constraint, while maintaining a model close to the real cluster, we consider that the platform is made of 486 homogeneous nodes of size $M = 128 \text{ GB}$.

The platform is managed by the SLURM scheduler, and as mentioned above, we observed the submission of I/O-intensive jobs. The cluster is often used for taxonomic identification of DNA fragments. To perform such identification, input databases¹ are often used, whose size vary from a few GB up to 800 GB. We also noticed the use of genome alignment tools that require to store large text files², and it is likely that they need to load multiple of them, resulting in high I/O demands. From this platform we extract historical workloads containing the following information.

B. Usage of real cluster logs

The logs we use contain historical data on jobs, namely their exact submission time, their requested walltime, their actual duration, the number of cores they required and the corresponding user's name. Since explicit data dependencies are not encoded in SLURM job specifications, we do not have access to the actual input files of these jobs. We thus create an artificial data dependency pattern that replicates user behaviors. Each job uses exactly one input file. We consider two jobs J_i and J_j . These jobs are considered to share their input file if they match the three following requirements:

¹See <https://benlangmead.github.io/aws-indexes/k2>

²See https://ftp.ensembl.org/pub/release-107/maf/ensembl-compara/multiple_alignments/10_primates.epo/

- 1) $Cores(J_i) = Cores(J_j)$, i.e., they request the same number of cores;
- 2) J_i and J_j are submitted by the same user;
- 3) J_i and J_j are submitted within an 800seconds time frame. We consider this timeframe to be a reasonable amount of time for a user to submit all of their jobs using the same input file.

Otherwise, we consider that J_i and J_j are using distinct input files. In theory, two users could share the same file. However, because they are using subsets of different databases, it is very unlikely that two users would work on the same project using the same databases, therefore we ignore this possibility.

We consider that these jobs are dedicated to processing their input file. Hence, the more cores the job requests, the larger its input file. This allows to estimate the size of files as follows:

$$Size(File(i)) = \frac{Cores(J_i)}{20} \times 128 GB$$

Indeed, if a user needs 128 GB of memory but would request only 1 of the 20 available cores, that user would block the node for all other users, for lack of remaining memory. Consequently, a user who needs the whole memory of a node will reserve all its cores.

The utilization levels in the log of the platform are typically high (> 90%), but not fully consistently so. The vast majority of jobs on these resources are single-node jobs and thus fit in our framework. The few multi-nodes jobs are not representative of the typical usage of the platform as they account for only 0.32% (6440 out of almost 2 million jobs) of the workload. They are replaced by as many single-node jobs as necessary to represent the same workload. We notice that jobs durations extend up to 10 days, while some jobs only last a few minutes. Even if the workload is not homogeneous, it is representative of the real usage from an actual user community including, but not exclusively consisting of, many subfields of the life sciences with highly data-dependent workflows.

In our experiments we evaluate our schedulers week by week. Our workload is constituted of 51 weeks, ranging from January the 3rd 2022 to December the 25th 2022. To target different scenarios but avoid scheduling the whole year, we randomly selected 12 weeks and extracted the jobs submitted within these weeks from the logs. Over these 12 weeks of workloads, we scheduled 1 986 496 jobs. The most important for a user is that all the jobs that have been submitted at once are finished as soon as possible. Thus we introduce the notion of *user session*. We define the stretch of a *user session* as the sum of the stretch of jobs submitted by a user in a 5-minute window, and we concentrate on this value when reporting results. There have been 1083 distinct users and 136 404 *user sessions*.

To simulate these jobs in a realistic, steady-state operation of the platform, we consider both jobs submitted before and after the week under consideration. More precisely, we proceed as illustrated in Figure 1.

Our evaluation is based on three metrics: the stretch of each *user session*, the total time spent waiting for a file to

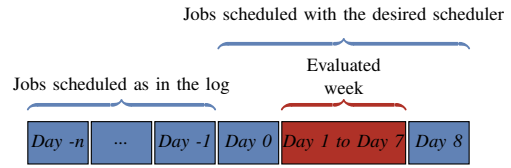


Fig. 1: Methodology followed to schedule and evaluate jobs from a specific week while avoiding edge effects.

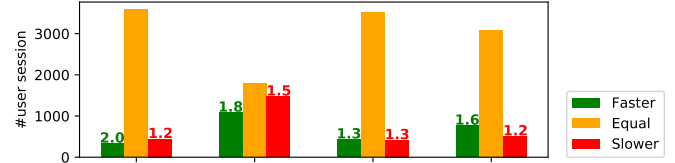


Fig. 2: User session stretch compared to FCFS on the workload of week 40 (low utilization case).

become available, as well as the total core time. We evaluate our proposed schedulers on these metrics and compare them to FCFS and EFT, with and without backfilling.

C. Simulator description

All strategies as well as the two baselines have been implemented on a simulator that we developed³. The simulator is an event-based simulator made in Python and focusing on data locality. As in real systems, the schedule is re-computed entirely at each unexpected event (job submission or job termination before the walltime). The scheduler is aware only of jobs that have been submitted before the current time and does not know the real duration of a job.

D. Results on an underutilized cluster

Among the 12 randomly-selected weeks, different workload conditions arise. We first concentrate on a low utilization case. Figure 2 shows the number of *user session* stretches that were faster, slower, or equal (+/- 1%) compared to FCFS. The number above the bars is the average speed-up (or slow-down) of the faster (or slower) *user sessions*. We observe that our schedulers mostly do not bring any improvement for most of the *user sessions*, except for LEA, which degrades the performance for more than a third of its *user sessions*.

To understand the issues LEA encounters on this workload, we need to study the cluster’s usage over time. Figure 3 shows the cluster usage when using FCFS. The vertical axis represents the number of requested cores either running on a node (lower half, the maximum is the total number of cores on our cluster: 9720) or in the queue of jobs waiting to be executed (top half, where the number of requested cores can be much higher than the cluster’s capability). The blue line shows the number of cores (used or requested) for all jobs. The red lines show the number of cores used by the evaluated

³Code, anonymized logs and the methodology followed to randomly select our evaluated weeks are available for reproducibility: <https://github.com/MaximeGonthier/Locality-aware-batch-scheduling>

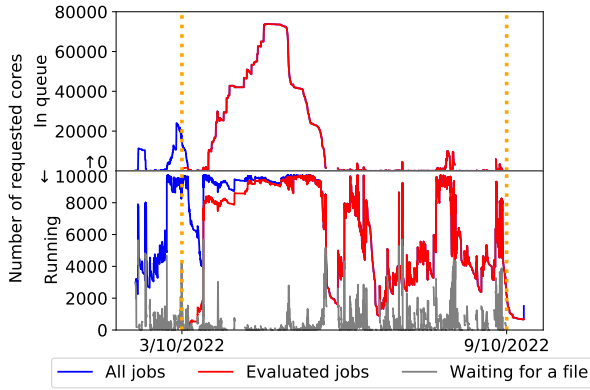


Fig. 3: Visualization of the utilization rate of the cluster on the workload of week 40 (low utilization case) with FCFS. For readability, the *in queue* and *running* scales are different.

Schedulers	EFT	LEA	LEO	LEM
Reduction from FCFS	0.0%	14.8%	0.8%	3.8%

TABLE I: Percentage of reduction in data transfer time relative to FCFS on the workload of week 40 (low utilization case).

jobs, i.e., those jobs that have their submission time within our evaluation window. Thus, this forms a subset of the set indicated by the blue line. If a line is present in the top half, it means that some jobs could not be scheduled and are thus waiting in the queue of available jobs. A node may have some available cores but not enough to accommodate some jobs with large requirements. This explains why the waiting job queue may not be empty even if the lower half does not reach the maximum. The gray line represents the number of cores currently loading a file. Lastly, the orange lines delimits the submission times of our evaluated jobs, in other words it is our evaluation window. By looking at the top half of Figure 3 we understand that the job queue is empty for more than half of the evaluation window. Consequently, FCFS and EFT are very efficient. The earliest available node is in most cases a node that can start the job immediately, explaining why most *user sessions* get no improvement in Figure 2. LEO is a strategy that uses the earliest available time t_k of a node to decide if it should compute a score like LEA, that puts a large weight on transfer time, or weigh equally t_k , the transfer and eviction durations. Thus, on underused clusters, LEO has a behavior close to EFT, while trying to minimize evictions. Similarly, LEM switches between EFT and LEA depending on the cluster’s usage. On this particular workload the cluster’s usage is under 100% more than half the time, so LEM behaves similarly to EFT. On the contrary, LEA favors data re-use over an early start time for a job.

Table I shows the reduction of the amount of time spent waiting for a file to be ready before starting the computation, relative to the total waiting time of FCFS. In this case, LEA is the only strategy that decreases this waiting time by re-using input files. It will create a queue of jobs that already

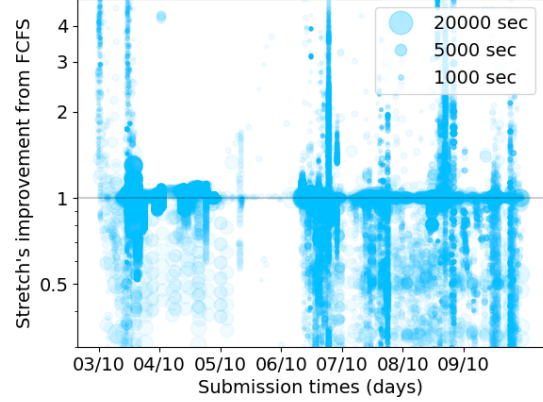


Fig. 4: Stretch times of each user session from week 40 (low utilization case) using LEA compared to FCFS. Cropped at 0.3 and 5. The max and min values are 543 and 0.003.

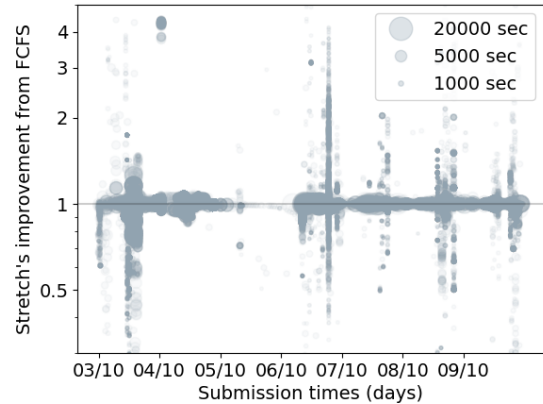


Fig. 5: Stretch times of each user session from week 40 (low utilization case) using LEM compared to FCFS. Cropped at 0.3 and 5. The max and min values are 543 and 0.007.

have a valid copy of their file loaded on a node, waiting to be able to re-use their input. On an underutilized cluster, this creates situations where some nodes are idle while submitted jobs wait in a queue, resulting in increased queue times. The consequences for the stretch are immediately distinguishable in Figure 4. This plot shows for each *user session*, the ratio of its stretch with LEA over its stretch with FCFS (hence a value above 1 means LEA improves the stretch). The size of a circle is proportional to the *user session*’s duration. On the workload of week 40, we observe several “columns” of jobs submitted at the same time (hence sharing the same file with large probability) with a ratio lower than 1, hence a worst performance with LEA. This means LEA is waiting to re-use the files before starting the jobs, whereas FCFS paid the cost of loading the file on other nodes, but started the jobs earlier than LEA, leading to shorter completion times. From Figure 5 concerning LEM, we see that most jobs have an

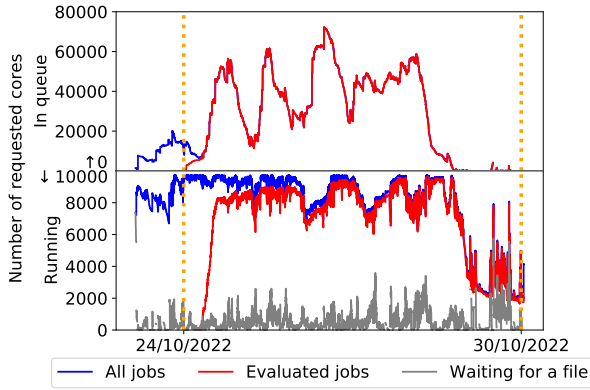


Fig. 6: Visualization of the utilization rate of the cluster on the workload of week 43 (high utilization case) with FCFS.

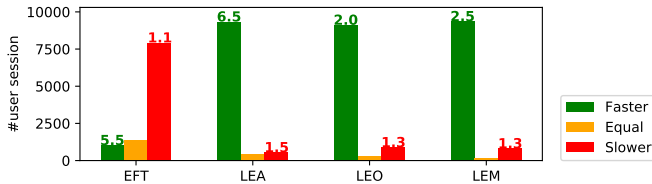


Fig. 7: User session stretch compared to FCFS on the workload of week 43 (high utilization case).

improvement close to 1, showing that LEM does not fall into LEA’s pathological case.

To summarize, on an underutilized cluster, LEA’s focus on locality does not allow optimal utilization of the cluster, while LEO and LEM, thanks to their flexibility, achieve performance close to EFT.

E. Results on a saturated cluster

We concentrate here on week 43, which we identify as a workload saturating the cluster. Indeed, Figure 6 shows a queue of several thousands requested cores for the whole duration of the evaluated week.

In this workload, LEA, LEO, and LEM reduce transfer time by 21%, 1%, and 11% compared to FCFS, respectively. Due to this reduction, we can observe in Figure 7 that the majority of *user session* stretches are improved by our proposed strategies

On a saturated cluster, filling all cores with the first jobs in the queue, as FCFS does, is not critical. It is more beneficial to group jobs that share the same file. This way, the first few jobs will have a longer queue time than with FCFS, but over time, file reuse creates a snowball effect that reduces the queue time of all subsequent jobs. Moreover, the queue contains enough jobs to fill all the nodes, even when grouped by input file. We thus avoid the pathological case outlined in Section V-D for LEA. LEA’s strategy, also found in LEO and LEM, allows to group jobs that share an input file, thus greatly reducing the stretch of each *user session*. LEA focuses only on locality, while LEO and LEM are mixed strategy, which explains the

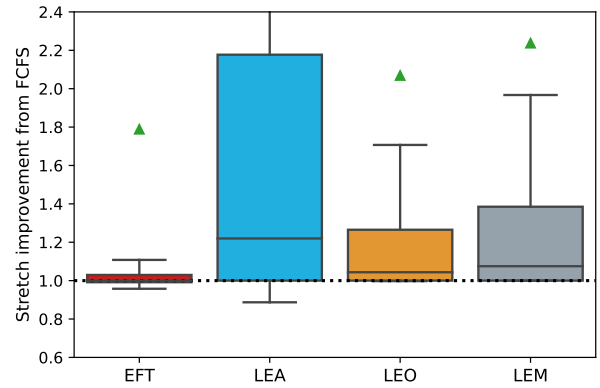


Fig. 8: Stretch’s improvement from FCFS on all evaluated weeks. The upper whisker of LEA extends up to 5.

Schedulers	EFT	LEA	LEO	LEM
Reduction from FCFS	0.9%	17.1%	0.9%	7.1%

TABLE II: Percentage of reduction in data transfer time relative to FCFS on all workloads.

higher average improvement for LEA compared to LEO and LEM in this very locality-favoring situation.

F. Complete results

We report here the aggregated results of the 5 schedulers over the 12 weeks evaluated, after scheduling nearly 2 million jobs. Figure 8 and 10 depict the distribution of the improvement of each of our *user session* stretch, compared to FCFS. In other words we represent $Stretch_{FCFS}/Stretch_{Scheduler}$ for each *user session* and each scheduler. The horizontal black dotted line corresponds to no improvements from FCFS, which is the situation where the sum of the queue time and transfer time of the *user session* is the same as with FCFS. The solid lines are the median and the triangles are the average improvements. Because outliers get a huge improvement, the averages are usually higher than the medians. Because of such extreme outliers, we are more interested in median values. An improvement above 1 is a speed-up. An improvement at 0.5 means that the sum of the queue times and transfer times for this *user session* was two times larger with the scheduler than with FCFS. For all boxplots shown in this paper, the box contains results within the [25%,75%] range, while whiskers are drawn at 12.5% and 87.5%.

a) *Results without backfilling*: We observe in Figure 8 that EFT does not bring any real improvement compared to FCFS. EFT takes into account file transfers when scheduling jobs but is a less aggressive strategy than LEA or LEM, and the savings in data transfers are only 0.9% (see Table II). EFT is not able to see that a large number of jobs using the same file should be scheduled on the same node: even if it would generate more queue time, the overall execution time would be lowered thanks to data re-use.

Schedulers	EFT	LEA	LEO	LEM
Reduction from FCFS	0.01%	0.37%	0.02%	0.10%

TABLE III: Percentage of reduction of the total core time used on all workloads relative to FCFS.

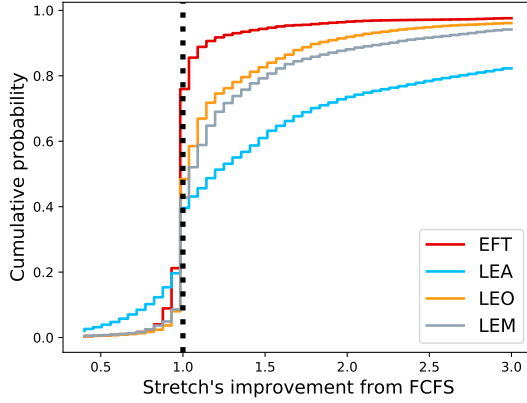


Fig. 9: Empirical distribution function of the stretch's improvement from FCFS on all evaluated weeks.

On the contrary, LEA has the largest median improvement. We can explain the larger median value for LEA from the good performance on heavily-saturated clusters (see Section V-E). Re-using the same files is not detrimental to the filling of all the nodes because there are enough jobs to cover all nodes. A large decrease (see Table II) in the time spent waiting for a file greatly reduces the stretch of each job.

LEM has a lower median. However, as can be seen in Figure 8, at least 87.5% of its results are above 1, i.e., an improvement, whereas for LEA, only approximately 75% of the results are above 1. LEM is a more versatile strategy and offers higher sustained performance on non-saturated cluster at the cost of fewer extreme improvements on heavily-saturated clusters.

Table III shows the reduction in terms of total core time (i.e. the core time on each core of each node over the 12 weeks). LEA reduces the total core time used by 49 000 core-hours over the 12 weeks (approximately 5h of the whole platform). Although it is a small percentage reduction, it can be of interest on large scale clusters that are highly demanding in terms of electrical resources.

From the same data shown in Figure 8, we plot an empirical distribution function in Figure 9. EFT's low variance is clearly visible in the sudden jump in probability around an improvement of 1. It is interesting to note that for LEA, 20% of the results are above an improvement of 300%. In addition, thanks to its change of scheduling strategy between EFT and LEA, LEM clearly reduces performance losses: it has a lower cumulative probability than LEA before the black dotted line.

Thus, we can conclude that LEA is the strategy that leads to more significant improvements: it is able to compute jobs between 1 and 2 times faster in 50% of the cases and between 2 and 6 times faster in 25%. It is slower in only 25% of

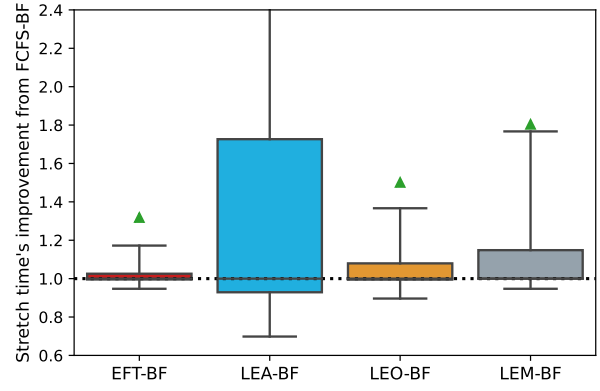


Fig. 10: Stretch's improvement from FCFS with backfilling on all evaluated weeks. The upper whisker of LEA extends up to 4.

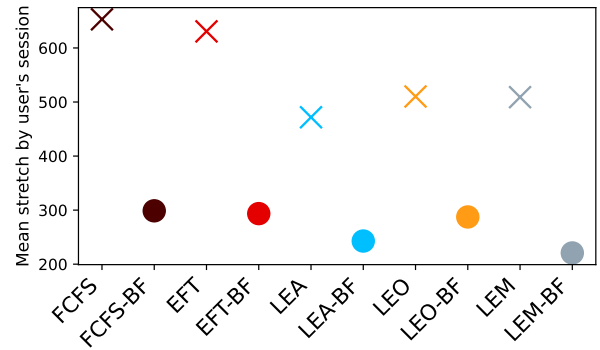


Fig. 11: Average stretch of all jobs evaluated on all 12 weeks by *user session* with (circles) and without (crosses) backfilling. Cropped at 200.

the workloads, and 12.5% of those are within a 0.15 slowdown. LEO is a more sustained strategy with 87.5% of its results with an improvement compared to FCFS, which shows that our opportunistic strategy is much more consistent, while still having great improvements in some cases. LEM is more versatile. It leads to an improvement in 87.5% of the results, with a speed-up of at least 7.5% in 50% of the results.

b) Results with backfilling: Figure 10 shows the results with the backfilling version of our schedulers and compared to FCFS with backfilling (FCFS-BF) on all workloads. We notice that our schedulers have smaller improvements with backfilling, because FCFS-BF already performs much better than FCFS. Figure 11 shows that with backfilling (circles), the performance difference compared to without backfilling (crosses) is much higher for FCFS-BF and EFT-BF. The reason is that our proposed strategies do not benefit from backfilling as much as FCFS and EFT for two reasons.

First, even if LEA, LEO and LEM consider data locality when backfilling, trying to fill a node as much as possible and optimizing data re-use are two contrary goals. Backfilling a job can compromise a re-use pattern that was planned by our locality-aware strategy, thus reducing the total amount of

re-used files.

Second, our strategies are already able to nicely fill the nodes without needing backfilling. Grouping jobs by input file implies that similar jobs end up on the same nodes. Jobs having the same duration and number of requested cores can much more easily fill a node to its fullest than a completely heterogeneous set of jobs. Consequently, FCFS-BF and EFT-BF already naturally benefit from increased data locality thanks to backfilling, leading to a reduced benefit in using LEA-BF, LEO-BF or LEM-BF.

Compared to FCFS-BF, our strategies still reduce the total queue time with backfilling but the difference is less significant.

Figure 10 shows that the improvement of EFT-BF compared to FCFS-BF is not significant. Out of our four heuristics, LEM-BF is the best compromise here. It is better than FCFS-BF in more than 75% of the cases, with 12.5% of those results above an improvement of 1.8. Among the slow-downs, only 12.5% are worse than 0.95.

VI. CONCLUSION

Batch schedulers are key components of computing clusters, and aim at improving resource utilization as well as decreasing jobs' response time. We have studied how one may improve their performance by taking job input file into account: in clusters dedicated to data analysis, users commonly submit dozens of jobs using the same multi-GB input file. Classical job schedulers are unaware of data locality and thus fail to re-use data on nodes. We have proposed three new locality-aware strategies, named LEA, LEO and LEM, capable of increasing data locality by grouping together jobs sharing inputs. The first one has a major focus on data locality, while the other two target a balance between data locality and load balancing. We have performed simulations on logs of an actual cluster. Our results show that LEA significantly improves the mean waiting time of a job, especially when the cluster is under a high computing demand. Without backfilling, LEA is better than our baseline in 75% of the cases (50% of the cases with backfilling). Our strategy called LEM is the best compromise. LEM is better than the baseline in more than 75% of the cases with or without backfilling (with a median improvement of 7.5% compared to our baseline without backfilling).

This work opens several exciting future directions. LEM can be tuned to better adapt to the utilization rate of the cluster. Switching between a locality-first and distribution-first strategy could be done gradually. An extension towards workflow scheduling where tasks depend on the output of a previous task would be very useful. Indeed the output file in memory would be re-used for other jobs, which is a new opportunity for locality-aware schedulers. In the long run, our objective is to consider other issues raised by batch scheduling like fairness between users, or dealing with advance reservations. Lastly, we want to integrate these strategies into real cluster schedulers to test their robustness in real-world situations. A first step would be to implement in the real system a labeling system that allows users to indicate the input data they are using. This

would allow for better input data for the simulation, which would strengthen the results of our simulation. We will then be able to use our techniques more confidently in place of the standard cluster schedulers.

REFERENCES

- [1] A. B. Yoo, M. A. Jette, and M. Grondona, "Slurm: Simple linux utility for resource management," in *Job Scheduling Strategies for Parallel Processing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2003.
- [2] N. Capit, G. Da Costa, Y. Georgiou, G. Huard, C. Martin, G. Mounie, P. Neyron, and O. Richard, "A batch scheduler with high level components," in *CCGrid 2005. IEEE International Symposium on Cluster Computing and the Grid, 2005.*, vol. 2, 2005, pp. 776–783 Vol. 2.
- [3] G. Staples, "Torque resource manager," in *Proceedings of the 2006 ACM/IEEE Conference on Supercomputing*, ser. SC '06. New York, NY, USA: ACM, 2006.
- [4] S. Kannan, M. Roberts, P. Mayes, D. Brelsford, and J. F. Skovira, "Workload management with loadleveler," *IBM Redbooks*, 2001.
- [5] R. L. Henderson, "Job scheduling under the portable batch system," in *Workshop on Job Scheduling Strategies for Parallel Processing*. Springer, 1995.
- [6] W. Gentsch, "Sun grid engine: towards creating a compute power grid," in *Proceedings First IEEE/ACM International Symposium on Cluster Computing and the Grid*, 2001.
- [7] Y. Etsion and D. Tsafir, "A short survey of commercial cluster batch schedulers," *School of Computer Science and Engineering, The Hebrew University of Jerusalem*, vol. 44221, 2005.
- [8] "Slurm workload manager," https://slurm.schedmd.com/sched_config.html, accessed: 2022-12-06.
- [9] D. Jackson, Q. Snell, and M. Clement, "Core algorithms of the maui scheduler," in *Job Scheduling Strategies for Parallel Processing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2001, pp. 87–102.
- [10] S. Srinivasan, R. Kettimuthu, V. Subramani, and P. Sadayappan, "Characterization of backfilling strategies for parallel job scheduling," 02 2002.
- [11] S. Leonenkov and S. Zhumatiy, "Introducing new backfill-based scheduler for slurm resource manager," *Procedia Computer Science*, 2015, 4th International Young Scientist Conference on Computational Science.
- [12] D. Jackson, Q. Snell, and M. Clement, "Core algorithms of the maui scheduler," in *Job Scheduling Strategies for Parallel Processing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2001.
- [13] D. G. Feitelson and M. A. Jette, "Improved utilization and responsiveness with gang scheduling," in *Job Scheduling Strategies for Parallel Processing*. Berlin, Heidelberg: Springer Berlin Heidelberg, 1997.
- [14] S. Parsa and R. Entezari-Maleki, "Rasa: a new grid task scheduling algorithm," *JDCTA*, vol. 3, 01 2009.
- [15] A. G. Delavar, M. Javanmard, M. B. Shabestari, and M. K. Talebi, "Rsdsc (reliable scheduling distributed in cloud computing)," *International Journal of Computer Science, Engineering and Applications*, 2012.
- [16] S. Ghanbari and M. Othman, "A priority based job scheduling algorithm in cloud computing," *Procedia Engineering*, vol. 50, no. 0, 2012.
- [17] J. Bent, D. Thain, A. C. Arpaci-Dusseau, R. H. Arpaci-Dusseau, and M. Livny, "Explicit control in the batch-aware distributed file system," in *NSDI*, vol. 4, 2004.
- [18] D. Borthakur *et al.*, "Hdfs architecture guide," *Hadoop apache project*, 2008.
- [19] J.-F. Weets, M. K. Kakhani, and A. Kumar, "Limitations and challenges of hdfs and mapreduce," in *2015 International Conference on Green Computing and Internet of Things (ICGCIoT)*, 2015.
- [20] P. Mishra, T. Agrawal, and P. Malakar, "Communication-aware job scheduling using slurm," in *49th International Conference on Parallel Processing-ICPP: Workshops*, 2020.
- [21] D. S. Nikolopoulos and C. D. Polychronopoulos, "Adaptive scheduling under memory constraints on non-dedicated computational farms," *Future Gener. Comput. Syst.*, vol. 19, 2003.
- [22] P. Agrawal, D. Kifer, and C. Olston, "Scheduling shared scans of large data files," *Proc. VLDB Endow.*, aug 2008.
- [23] M. A. Bender, S. Chakrabarti, and S. Muthukrishnan, "Flow and stretch metrics for scheduling continuous job streams," in *Proceedings of the Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, ser. SODA '98. USA: Society for Industrial and Applied Mathematics, 1998, p. 270–279.