

Improving Batch Schedulers with Node Stealing for Failed Jobs

Yishu Du^{1,2} | Loris Marchal¹ | Guillaume Pallez³ | Yves Robert^{1,4}

¹Laboratoire LIP, ENS Lyon, France

²Tongji University, Shanghai, China

³Inria Rennes, France

⁴University of Tennessee Knoxville, TN, USA

Correspondence

Yves Robert

Email: yves.robert@ens-lyon.fr

Funding information

Yishu Du is supported by grants from the Chinese Scholarship Council and from Inria.

After a machine failure, batch schedulers typically re-schedule the job that failed with a high priority. This is fair for the failed job but still requires that job to re-enter the submission queue and to wait for enough resources to become available. The waiting time can be very long when the job is large and the platform highly loaded, as is the case with typical HPC platforms. We propose another strategy: when a job J fails, if no platform node is available, we steal one node from another job J' , and use it to continue the execution of J despite the failure. In this work, we give a detailed assessment of this node stealing strategy using traces from the Mira supercomputer at Argonne National Laboratory. The main conclusion is that node stealing improves the utilization of the platform and dramatically reduces the flow of large jobs, at the price of slightly increasing the flow of small jobs.

KEYWORDS

Batch scheduler, parallel jobs, failures, resilient schedule, node stealing, maximum flow, mean flow.

1 | INTRODUCTION

Batch schedulers, a.k.a Resource and Job Management Systems (RJMS), are a key component of the supercomputing infrastructure. Users make reservation for their parallel job that includes information such as an upper-bound on the

expected length (called the wall time), and the desired number of resources needed for the execution. The task of the batch scheduler is then to allocate these jobs on the computing platform, with the end goal of optimizing some metric or combination of metrics.

In the last decade, batch schedulers have faced additional constraints: on state-of-the-art platforms, an increasing number of users experience the crash of a node belonging to their reservation set during the execution of their job. This is because platforms are composed of more and more nodes to accommodate for an endless increase in job demands. This scaling is the main reason for the increasing number of failures. The most recent supercomputers such as Frontier, Fugaku, or LUMI (the top three entries in the TOP500 ranking [29]) are now embedding millions of cores (with a peak at 10.6M for Sunway TaihuLight (ranked 6th)). These colossal systems are prone to failures: even if each of their cores has a low probability of failure, the failure probability of the whole system is much higher. More precisely, assume that the *Mean Time Between Failure* (MTBF) of each computing resource is around 10 years, which means that such a resource should experience an error only every ten years on average, and which shows that computing resources are quite reliable individually. When running a simulation code on 100,000 of these resources in parallel, the MTBF is reduced to only 50 minutes [13]: on average one of the resource crashes every 50 minutes. With one million of such resources, the MTBF gets as small as five minutes, while codes deployed on such extreme-scale platforms usually last for hours or days. As the demand for computing power increases, failures cannot be ignored anymore, and fault-tolerant mechanisms must be deployed, such as checkpoint/restart mechanism.

When a job fails, the standard policy is to relaunch it as soon as possible (from its last valid checkpoint): the job is put back in the submission queue, but with a high priority, so that it can be re-executed rapidly (e.g., see the 'job failover' section in [15]). If there is a free node available at the time of the failure, the failed job will be able to resume execution (almost) immediately: because it is given a high priority, the failed job will be re-assigned all the surviving nodes of its reservation, plus the free node. Of course it may well be the case that no free node is available at the time of a failure, say if the platform is over-subscribed. In that case, the failed job will have to wait until enough resources become available for its re-execution.

In this paper, we propose a novel approach for High Performance Computing (HPC) platforms: if there is no free node available when a failure strikes a job, we propose to create one! This means to interrupt another job that is currently executing, and to *steal* one of its nodes and assign it as a new resource to the failed job. This *node stealing* approach is inspired by similar ideas in cloud computing, where users who have paid for *spot instances* [21, 17, 19] can have their resources taken back without prior notification. To the best of our knowledge, this work is the first work that studies node stealing in an HPC framework. There are several decisions to explore:

- **Which job to interrupt?** Clearly, small jobs with one or few nodes are good candidates, because they are easier to re-schedule. But interrupting a small job whose waiting time is already high may not be fair to the owner of that job, so trade-offs between different optimization metrics must be achieved.
- **When to interrupt?** Immediately after the failure is the simplest solution, but the interrupted job will lose the work done since its last checkpoint. Another solution is to wait for a checkpoint before the interruption, or immediately enforce a proactive checkpoint, depending upon what is possible.

The main contributions of this paper are the following:

- A thorough description of the problem, and how to measure its usefulness;
- A focus on SFSJ (*Steal From Small Jobs*), a strategy which chooses the job to interrupt among those with the smallest number of nodes and, if ties, with the shortest execution time so far;

- An evaluation of SFSJ in a simulated framework, based upon trace-based scenarios;

Furthermore, we provide elements about a comparative assessment of several other node stealing strategies (the complete data is not included due to lack of space).

The rest of the paper is organized as follows. Section 2 provides motivational data on the impact of failures while Section 3 works out a toy example. Section 4 details the design of the SFSJ strategy. Sections 5 and 6 are devoted to a comprehensive experimental comparison of SFSJ: Section 5 presents the methodology and the various potential objectives, while Section 6 presents the results and assesses the efficiency and limits of the approach. Section 7 discusses several other node stealing strategies. Section 8 surveys related work. Finally, Section 9 gives concluding remarks and hints for future work.

2 | MOTIVATION

This section provides a brief motivation for node stealing techniques in the presence of failures, which dramatically increase the flow of large jobs. Recall that in the scheduling literature [2, 3], the flow (or flow-time) of a job represents the total time spent by the job in the system, from submission to completion, and includes both waiting time and execution time¹.

We have simulated an execution of the workload submitted to the Mira platform at Argonne National Laboratory [6, 24, 22] in June 2017 (see details of the simulation in Section 5). Figure 1 shows different job flows for this execution. The x-axis corresponds to the *job size*: jobs are classified in categories depending on their requested number of nodes. Figure 1 (left) shows the maximum flow as a function of a job size, i.e., the maximum flow observed for jobs of a given size. Figure 1 (right) shows the mean flow as a function of a job size, i.e., the average flow observed for jobs of a given size. The central value of the boxplot represents the median, while the box extends from the lower to the upper quartiles. The upper whisker extends from the hinge to the largest value no further than $1.5 \times \text{IQR}$ from the hinge, where IQR is the inter-quartile range or distance between the first and third quartiles. The lower whisker extends from the hinge to the smallest value at most $1.5 \times \text{IQR}$ of the hinge. Data beyond the end of the whiskers are called "outlying" points and are plotted individually. Box plots provide flows over five randomly generated failure scenarios.

On both subfigures, a red dot corresponds to the flow obtained in a failure-free environment. We can observe that larger jobs have larger flows in a failure-free environment.

We have enriched the figure with the flows of the same jobs in presence of failures, assuming that the MTBF of the platform is one hour (which is the typical MTBF expected for future scale systems [23]). Given that Mira platform had 49152 nodes, this leads to an individual MTBF for each node of $\mu_{\text{ind}} = 5.61$ years. Failures are randomly generated following a Poisson process (parameter $\lambda = 1/\text{MTBF}$) and several failure scenarios are considered. The results are reported in (black) box plots. Jobs are checkpointed according to the optimal Young/Daly period $P_{YD} = \sqrt{2 \frac{\mu_{\text{ind}}}{p} C}$, where p is the job size and $C = 5$ minutes is the (assumed) checkpoint length for all jobs. When a job experiences a failure, it is re-scheduled using the baseline strategy: the job is put back into the queue with highest priority, meaning that it will be re-executed as soon as enough nodes (the job size) are available. If there is a free node available at the time of the failure, this free node can 'replace' the node struck by the failure, and the failed job will be able to resume execution almost immediately. Of course this leads to re-scheduling all the jobs in the execution queue that have not yet started their execution. This baseline strategy is the one used in several batch schedulers such as IBM's LSF [15].

¹In the operating system literature, another word for job flow is *turnaround time* [25].

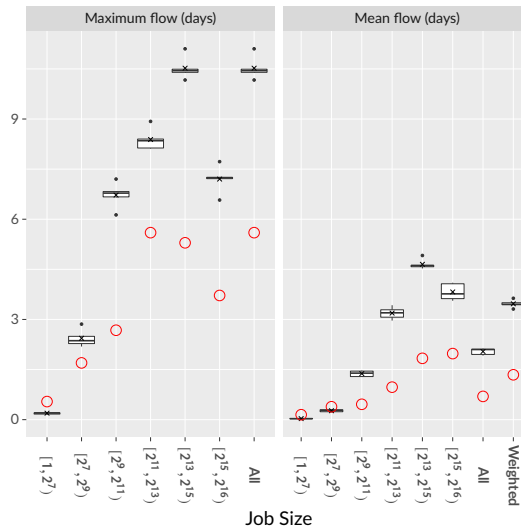


FIGURE 1 Maximum flow and mean flow as a function of job size, without failures (red dots) and with failures (box plots), using Baseline (workload: Mira, June 2017 [22]). “Weighted” mean flow uses job sizes as weights.

Several observations from Figure 1 can be made:

1. The impact of failures is dramatically higher for jobs with more than 2048 nodes, whose flow has increased much more than the flow of jobs with less than 512 nodes. This is because large jobs are harder to re-schedule, due to their high resource demand.
2. The flow of short jobs may be reduced by failures. Indeed, when a large job fails, it has to wait for the completion of another job to get a spare node. During this waiting time, many nodes are left idle and can be used by small and short jobs using backfilling: short jobs are allowed in the “holes” of the schedule; they may start earlier than some (longer) jobs submitted before them, provided that they do not delay these jobs.

This observation that failures have a non-uniform impact on jobs of different sizes, is at the heart of our approach: would it possible to *steal* nodes from small jobs when large jobs are struck by failures, in order to mitigate the increase of large job flows? A key contribution of this work is to assess the efficiency of *node stealing* in various execution scenarios. Intuitively, if the platform is not over-subscribed, idle nodes will be available most of the time, and node stealing will be rarely (if at all) needed. But as the subscription rate augments, we expect node stealing to become more frequent.

3 | TOY EXAMPLE

This section works out a toy example to explain how node stealing can be used to decrease these large job flows. It provides insight in the decision process made throughout this paper. Consider a platform with 8 nodes. Five jobs are released at time $t = 0$: see Table 1 and Figure 2 for details on these jobs. Since all five jobs are released simultaneously at time $t = 0$, we can assume that the scheduler has broken ties so that the jobs are scheduled in the order J_1, J_2, \dots ,

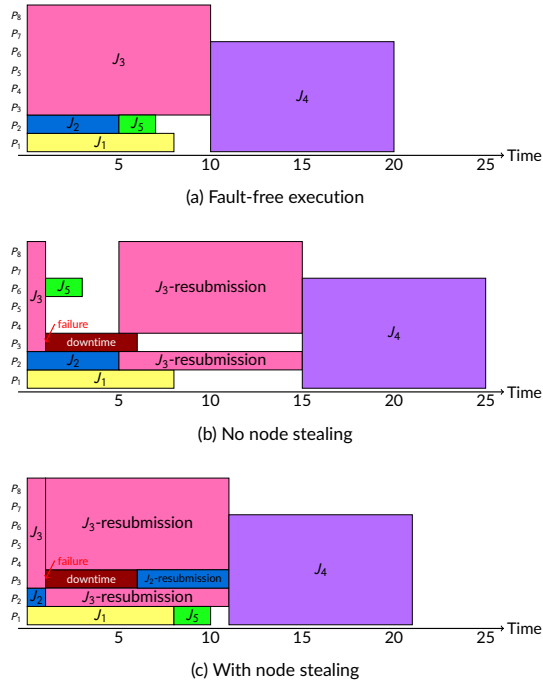


FIGURE 2 Toy example, job details in Table 1. Subfigures (b) and (c) assume that a failure occurred at $t = 1$ on P_3 .

up to J_5 .

At time $t = 0$, the scheduler starts J_1 on P_1 , J_2 on P_2 , and J_3 on P_3 to P_8 . It reserves P_1 to P_6 for J_4 at $t = 10$. At time $t = 5$, it backfills J_5 on P_2 since it will not delay J_4 . Figure 2(a) depicts the fault-free execution.

We consider now that the platform will experience failures. To simplify the example, jobs are not checkpointed and can resume immediately after a failure if there are available nodes, meaning that we neglect any recovery cost. Downtime (rejuvenation time) for each node is $D = 5$, meaning that a node struck by a failure at time t is up again at time $t + 5$. Suppose then that a failure strikes P_3 at $t = 1$. Figure 2(b) depicts the standard scenario. J_3 fails at $t = 1$ and is now the job with the highest priority for re-scheduling. There are only five free nodes at $t = 1$, and this holds true until $t = 5$. Hence J_3 is scheduled for execution at $t = 5$ on nodes P_2 and P_4 to P_8 (since P_3 is unavailable until $t = 6$ due to downtime). Now J_3 completes at time 15 and J_4 completes at time 25. Using backfilling, J_5 is scheduled at $t = 1$ on one available node (P_6 in the figure). In line with the observations made in Section 2, we see that the smallest job has finished earlier in the presence of a failure than without one, while the large jobs have suffered the most from the failure.

What happens instead if we steal a node when the failure strikes P_3 at $t = 1$? We represent this new scenario in Figure 2(c). At $t = 1$, we steal P_2 and thereby interrupt job J_2 . Job J_3 can re-execute immediately on nodes P_2 (replacing P_3) and P_4 to P_8 . J_3 now finishes at time 11. Then J_2 has highest priority and can re-execute on P_3 when it is up again at time 6. Now J_2 completes at $t = 11$. Then J_4 is scheduled at time 11 and completes at time 21. Using backfilling, J_5 executes on P_1 when it becomes available.

Table 1 reports some statistics about the flows of the five jobs in the different scenarios without or with node stealing. We see that the flows of the large jobs J_3 and J_4 have decreased, at the price of increasing the flow of

TABLE 1 Job information for the toy example.

id	Release time	Job size	Job length	Flow without node stealing	Flow with node stealing
J_1	0	1	8	8	8
J_2	0	1	5	5	11
J_3	0	6	10	15	11
J_4	0	6	10	25	21
J_5	0	1	2	3	10

the small jobs J_2 and J_5 . The maximum value of the flow has decreased from 25 to 21. However its mean value has increased from 11.2 to 12.2. This is interesting as it shows that the mean flow is highly influenced by small jobs, while these jobs are not the most critical jobs on HPC platforms. Another widely used metric is the weighted mean flow, where the mean is weighted by the number of nodes of each job. Here, the weighted mean flow without node stealing is $(1 \times 8 + 1 \times 5 + 6 \times 15 + 6 \times 25 + 1 \times 3) / (1 + 1 + 6 + 6 + 1) = 17$, while the one with node stealing is 14.733.

We also see that the total idle time of the 8 nodes has decreased. Altogether, node stealing seems quite beneficial here! Beyond this toy example, a major contribution of this paper is to assess the usefulness of node stealing in various realistic execution scenarios.

4 | NODE STEALING

This section provides a high-level description of the classic conservative backfilling strategy used by batch schedulers (Section 4.1) and details how to extend it to implement node stealing (Section 4.2).

4.1 | Baseline strategy

First-Come First-Serve (FCFS) is a simple approach to submit jobs on parallel supercomputers. However, FCFS often leads to a waste of resources: when there are not enough free nodes for the next job, these free nodes remain waiting until additional nodes become available. A widely-used solution is to use non-FCFS policies, i.e., to allow for a (limited) reordering of the jobs in the queue. *Backfilling* schedulers [20] have been proposed to allow small jobs further away in the queue of waiting jobs to be processed whenever there are enough resources for them. Backfilling may lead to delay some previously allocated jobs, hence it must be controlled so as to guarantee that large jobs will get processed eventually. This is why, in the *conservative backfilling* algorithm, short jobs are moved ahead only if they do not delay any previous job already scheduled.

When a failure hits the system, the remaining part of the job that failed is put back into the scheduling queue, with the highest priority. Depending upon the absence or presence of a resilience mechanism, the remaining part of the job can represent either the whole job or the fraction of the job after the last checkpoint. The schedule is then recomputed with all jobs that have not started yet. If there are multiple jobs that have failed in the queue, they are sorted by non-decreasing arrival time. Throughout the paper, Baseline will denote this conservative backfilling scheduling strategy.

4.2 | Node stealing protocol

Node stealing should be seen as a feature that can be added on top of any batch scheduling strategy. In this work, we add this feature on top of Baseline scheduling. The core idea is the following: when a failure hits a job (say job J_1), and if there is no (free) node available at the time of a failure, then we select another job (say job J_2) which we *interrupt*. A node from job J_2 is allocated to job J_1 , so that job J_1 can resume its execution immediately, either from its last checkpoint if any, or from scratch. Job J_2 is then marked as failed, and it is restarted, again from its last checkpoint if any, otherwise from scratch. The schedule is then recomputed with the following priorities: (high) job J_1 ; (medium) job J_2 ; (low) other submitted jobs in the order of the underlying scheduling algorithm (here Baseline).

In the following sections, we focus on a single node stealing strategy and select the job to interrupt (called *victim* in the following) using the following procedure: *among all running jobs that use the fewest nodes, we select the one that has been submitted the latest*. In other words, the selection criteria are job size first, and job release time to break ties. If no victim job is found with fewer nodes than the failed job, node stealing is not activated. Throughout the paper, we let SFSJ (*Steal From Small Jobs*) denote this particular node stealing strategy. Other node stealing strategies are discussed and evaluated in Section 7, along with the possibility to take a proactive action, i.e., checkpoint the job chosen to be interrupted before actually interrupting it.

We point out that Baseline and SFSJ behave exactly the same when a free node is available when a job is struck by a failure. Both strategies have the failed job re-submitted with high priority, and therefore start re-execution immediately. However, when no free node is available when a job is struck by a failure, the strategies differ: Baseline lets the the failed job wait until enough resources become available, while SFSJ interrupts another job to be able to restart the failed job immediately. Besides, SFSJ has the same complexity than the Baseline strategy, which $O(n)$ to schedule n jobs.

5 | EVALUATION METHODOLOGY

In this section, we detail the evaluation setup. Our approach relies on the Batsim simulator [10], which emulates a batch scheduler on a parallel platform (see Section 5.1). We have extended Batsim to simulate a failure-prone environment. This extension uses a platform size and a job trace as input. We emulate the Mira supercomputer at Argonne National Laboratory using public traces of this machine [22]. The details of the traces and how they are modified to incorporate resilience mechanisms are presented in Section 5.2. Finally, we discuss key objectives used to evaluate the performance of batch schedulers in Section 5.3.

5.1 | Simulation environment

Our simulation environment relies on the existing Batsim simulator and the Batsched scheduling algorithm toolbox. Batsim (Batch scheduler SIMulator) [10] is a modular RJMS simulator based on SimGrid [5], which is a state-of-the-art distributed platform simulator. Batsim is in charge of simulating the behavior of the computational resources. Batsched is a C++ toolbox of scheduling algorithms that take decisions on when and where (which nodes) to schedule a job, and possibly when to interrupt a job. Batsched communicates with Batsim to receive the information about released jobs and to send scheduling decisions.

There already exists an event injection mechanism in Batsim/Batsched that allows to make the scheduler aware of external events on the platform. We used and adapted this mechanism in order to simulate node failures and rejuvenation. Whenever Batsched receives the message that a given node has failed, this node is removed from the

TABLE 2 Number of jobs categorized by size (requested number of nodes).

Job size intervals	1	$[2^1, 2^3]$	$[2^3, 2^5]$	$[2^5, 2^7]$	$[2^7, 2^9]$	$[2^9, 2^{11}]$	$[2^{11}, 2^{13}]$	$[2^{13}, 2^{15}]$	2^{15}	49152	total
(Mira) June 2017	8	2	6	10	74	2103	809	269	22	8	3311
(Mira) March 2018	31	3	6	69	117	2481	923	350	31	13	4024
(Intrepid) June 2013	0	0	0	0	0	2001	574	362	31	2	2970

set of machines available for computations, and thus cannot be used for executing jobs. If a job was running on the node that just failed, Batsched notifies Batsim that this job is interrupted. Besides, the whole schedule predicted by Batsched has to be recomputed from the current time (the failure time). In the new schedule, the remaining fraction of the job interrupted is given a higher priority, as detailed in Section 4.2.

Similarly, when Batsched receives the message that a node that failed before has been rejuvenated, it adds this node to the set of available machines, and the whole schedule is recomputed to take advantage of this newly available resource. Note that in steady-state mode, not all nodes will be up: some have been struck by failures and are rejuvenating. Hence, huge jobs are likely to wait for longer times before execution, and some may fail and be re-submitted several times.

The code developed to run these simulations is publicly available [9], together with Python scripts used to generate failures and R scripts used to handle workloads traces, analyze the results and produce the final plots.

5.2 | Supercomputer workloads

We use traces of the Mira and Intrepid supercomputers [22] to evaluate the performance of node stealing. Specifically, Batsim uses the following data to compute its schedule: (i) release time; (ii) wall time (predicted execution time of the jobs); (iii) length (actual execution time of the jobs, also called *delay* by Batsim); (iv) number of nodes. We conduct the experiments on two months from the Mira trace: June 2017 and March 2018. These months were selected because their *stress*² on the platform are quite reasonable (89.63% and 97.78% respectively), as well as sufficiently different to represent different usage scenarios. Job sizes for both months are detailed in Table 2. We also performed experiments on the June 2013 month from the Intrepid trace. This month is selected as it does not contain any job with less than 2^9 nodes, hence allowing to test our strategy on a quite different workload. Its *stress* on the platform is 89.55%.

In order to evaluate the impact of failures, we had to transform the traces and control the fault-tolerance mechanism. The full script that takes as input a trace and returns the modified trace is publicly available [9]. The first step is related to the incorporation of failures. Given that Mira platform has 49152 nodes (resp. 40960 for Intrepid), and because we consider failure-intensive scenarios where one or several resources can be down at any time, we reduce the size of the largest jobs from 49152 nodes to 49000 nodes (resp. from 40960 to 40900 for Intrepid). This ensures that no job is rejected because it requires more nodes than actually available on the platform. Finally, we assume that no job is interactive in the traces, for the following two reasons: (i) we cannot distinguish interactive jobs from other jobs in the traces; and (ii) the scheduler would typically exclude interactive jobs from the set of jobs that should be considered in the node-stealing approach.

We can measure the utilization of the platform in a failure-free scenario for this new workload using Baseline. Unsurprisingly, the utilization is lower than the stress, and notably for March 2018, because of scheduling constraints. We present this data along with statistics about job length in Table 3.

The second step is to add fault-tolerance mechanisms to job submission data. Failures are randomly generated

²The stress is defined as the sum of the lengths of jobs submitted this month divided by the total platform availability time in the month, including all nodes.

TABLE 3 Workload utilization and job lengths.

	Failure-free utilization	Number of jobs	job execution time			
			min	max	mean	median
(Mira) June 2017	88.51%	3311	55s	49.88h	2.64h	0.88h
(Mira) March 2018	92.88%	4024	26s	24.02h	2.79h	1.07h
(Intrepid) June 2013	89.90%	2970	26s	47.08h	2.55h	0.42h

following a Poisson process on each node with parameter λ_{ind} . The Mean Time Between Failure (MTBF) of each individual node is $\mu_{\text{ind}} = \frac{1}{\lambda_{\text{ind}}}$. The MTBF of the whole platform is $\mu = \frac{\mu_{\text{ind}}}{N}$ [13], where N is the total number of nodes of the platform. We assume that the system performs periodic checkpointing using the Young/Daly formula [31, 7]. This means that each job performs a checkpoint every $P_{YD} = \sqrt{2\mu_{\text{job}}C}$ units of time, where C is the time to perform a checkpoint (we use $C = 5$ minutes for all jobs), and μ_{job} is the MTBF for this job. Here μ_{job} is job dependent as it relies on the number of nodes p used by the job: we have $\mu_{\text{job}} = \mu_{\text{ind}}/p$ [13]. Given a periodic checkpoint strategy, the number of checkpoints to be taken linearly depends upon the length of the job. Hence we increase the length of each job accordingly. Furthermore, from a platform perspective, it is only natural to increase the wall time t_{walltime} in a similar way. We compute the new job execution time $t_{\text{exec}}^{\text{ckpt}}$ and new wall time $t_{\text{walltime}}^{\text{ckpt}}$:

$$t_{\text{exec}}^{\text{ckpt}} = t_{\text{exec}} + \left\lceil \frac{t_{\text{exec}}}{P_{YD}} \right\rceil \times C \quad (1)$$

$$t_{\text{walltime}}^{\text{ckpt}} = t_{\text{walltime}} + \left\lceil \frac{t_{\text{walltime}}}{P_{YD}} \right\rceil \times C \quad (2)$$

During execution, when a failure occurs, jobs are restarted from their last successful checkpoint.

Two key parameters to assess the performance of Baseline and SFSJ are the downtime and the platform MTBF. We conduct a detailed analysis of the impact of these parameters in Section 6.

5.3 | Measuring performance

When considering the performance of a batch scheduler, there are several metrics to assess. As already stated, from the user's perspective, the most important metric is to minimize the flow, or response time, of the job: "How fast can I get my results?". The flow is defined as the time elapsed from the initial submission of the job up to its completion, possibly after some unsuccessful attempts due to failures. However, from the platform owner's perspective, the most important metric is to maximize the utilization: "How much work can be executed on the platform per time unit?". The utilization is loosely defined as the fraction of time where nodes are doing useful work, i.e., make actual progress in the execution of some jobs. In the following, we provide more details on both metrics and detail how we modified the trace to provide a fair evaluation.

5.3.1 | Maximum and mean flow

The flow of a job represents the time spent by the job in the system. The flow is composed of two elements that add up, the waiting time (time elapsed from its submission to the start of its execution) and the execution time (time spent computing with the reserved nodes). If the job fails during execution, it is resubmitted to the scheduler, which usually gives a high priority to re-execution. If the job is checkpointed, only the remaining part of the job after the last checkpoint will be re-executed. Regardless, the job flow accounts for all re-executions and is computed from

submission until complete (successful) execution.

Usually, the user makes a reservation with a duration (called *wall time*) and a node count; it is their responsibility to ensure that the reservation has longer duration than the (expected) execution time. This may lead to over-length reservations, in particular when the user is only billed for execution time, not reservation time – a standard scenario on today's platforms. However, longer reservations usually experience a longer waiting time, which is an incentive for users to accurately estimate their reservation length.

Maximum flow is the largest flow for any job running in the system. Mean flow is the average over all jobs in the system. The weighted mean flow is the weighted average over all jobs, where each job is weighted by its size (its number of nodes). This latter quantity gives a higher weight to jobs that use a large number of nodes, which are typically the target jobs deployed onto supercomputers.

5.3.2 | Utilization

The utilization is defined as the ratio of the core-hours *occupied to progress a job* over the core-hours *available* during that period. One could expect an utilization close to 1 on a highly-subscribed platform. However, the two main factors that decrease utilization are the following:

1. Idleness due to scheduling: even with sophisticated backfilling techniques, large jobs bring specific constraints to the scheduler; not all nodes can be used at every instant.
2. Failure mitigation: the time spent to checkpoint jobs, to recover from a failure, and to re-execute fractions of jobs that have been lost (after the last checkpoint up to the failure) all decrease platform utilization. In addition, the time spent to re-execute fractions of jobs that have been lost (after the last checkpoint up to the failure) also decreases utilization. It is important to exclude failure mitigation techniques (such as checkpointing) from the utilization of the platform. Otherwise, an artificial way to increase the utilization would be checkpoint extremely often, hence reducing the waste after each failure.

While idleness due to scheduling has been studied for decades, failure mitigation is a more recent concern. Checkpointing jobs using the Young-Daly formula minimizes the overhead due to failure mitigation. However, resubmitting failed jobs induces an extra burden on the scheduler.

5.3.3 | Pruning the traces

Since we simulate a given month of the traces of the Mira platform, the platform is not fully loaded at the beginning of the simulation (first days of the month), and the values for utilization and flow of the jobs that completes are not representative. Similarly, as job submissions stop at the end of the month, the results (utilization and job completion times) are not meaningful after the last submission. Hence we have to carefully select the data used to compute appropriate statistics.

To compute the utilization of the platform as well as the fraction of time spent in various operations (computing, checkpointing, etc.), we define a time window, going from the 11th day of the month up to the 30th of the month when all activities are registered.

When measuring job flow, we cannot use the same time window: by considering only jobs that complete in a predetermined time window, we would not measure the performance of the same subset of jobs for different heuristics. We thus select a slightly different set of data: we order jobs by submission time and remove the first 20% of jobs

(intuitively, the ones that are submitted before the platform is fully utilized) as well as the last 20% of jobs (intuitively, the ones that completes later than the last submission time) and compute the flows of all remaining jobs.

6 | RESULTS

In this section, we describe the experiments that compare node stealing with the baseline strategy (conservative backfilling). As already stated, we perform simulations on the Mira workloads in June 2017 and March 2018, and the Intrepid workloads in June 2013 [22]. In Subsection 6.1, we start by demonstrating the usefulness of the node stealing approach in one specific scenario, for which both the MTBF and the downtime are equal to one hour. This allows us to qualitatively discuss the impact of the strategy. Then we move to a more thorough and quantitative evaluation with varying MTBF and downtime values in Subsection 6.2.

6.1 | Baseline Scenario: MTBF=downtime=1 hour

We start this scenario by a remark on the checkpoint of small jobs. With a MTBF of 1h (i.e. $\mu_{\text{ind}} = 5.61$ years), and a checkpoint size of 5 minutes, the Young/Daly period for a job running on 128 nodes is $\sqrt{2 \cdot \frac{5.61 \cdot 365 \cdot 24 \cdot 60}{128} \cdot \frac{5}{60}} = 8$ hours. This means that any job running on (or on less than) 128 nodes, and which lasts less than 8 hours never checkpoints. In practice, in the timeframe that we are studying, no job of less than 128 nodes performs any checkpoint, and less than 5% of jobs with [128, 512) nodes perform at least a checkpoint. For this scenario, we first discuss platform utilization and then flows.

6.1.1 | Utilization

TABLE 4 Platform useful utilization for various workloads.

Heuristic	(Mira) 06/2017	(Mira) 03/2018	(Intrepid) 06/2013
Baseline	79.32%	77.39%	76.71%
SFSJ	80.69%	78.50%	77.26%

The utilization is presented in Table 4. With SFSJ, it is 1.4 to 1.7% higher than with baseline scheduling, which is a positive gain, yet limited. To better understand this observation, in Figure 3, we report the fraction of total platform time spent into something else than “useful” computations: idle time, resilience mechanisms (checkpoints and restart), downtime, work wasted due to failures (un-checkpointed work when a failure strikes), and any waste due to node stealing (un-checkpointed work interrupted by node stealing and additional recovery time for applications killed).

Utilization gain can only come from a reduction of the idle time. For Baseline it corresponds to 5.2% of the platform usage for June 2017 in Mira, (5.8% for March 2018 in Mira and 7.3% for June 13 in Intrepid). Figure 3 corroborates the small utilization gains, however it shows that they correspond to relatively important reduction of platform idle time (20% in Mira-March 2018, 40% in Mira-June 2017 and 10% in Intrepid-June 2013). This first item shows that SFSJ is quite impactful given its leeway, especially when the workload contains small jobs.

We further observe that the additional overhead due to SFSJ (work wasted due to job interruption and additional recovery times) is negligible (around 0.1% for all months, as shown by the thin black line on Figure 3). This shows that additional resilience mechanisms that one could envision for node stealing (such as proactive checkpoint before

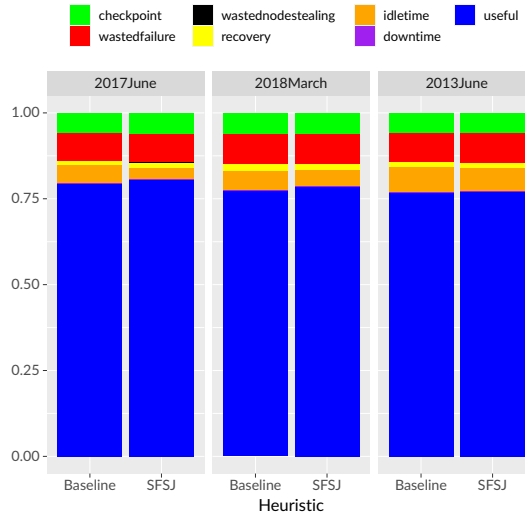


FIGURE 3 Decomposition of platform usage with Baseline and SFSJ.

interruption) have little room for improvement (see Section 7 for a discussion of other node stealing variants).

How often node stealing is actually used? Node stealing is only used when there is no free node available at the time of a failure. Table 5 provides some key statistics averaged over five randomly generated failure scenarios. Table 5a reports the percentage of time at least one free node is available right after a failure for both approaches. As shown in Table 5b, there is actually a free node available right after a failure, for 84% of failures in June 2017 (Mira), 89% in March 2018 (Mira) and 90% in June 2013 (Intrepid). In this vast majority of cases, node stealing is not activated, and both Baseline and SFSJ will resubmit the failed job with high priority, hence start its re-execution (almost) immediately. Finally, the different percentages for the reduction of idle time (Figure 3) can be explained by the different percentage of situations where SFSJ has to interrupt a job.

To conclude from a system performance perspective, there is only little room for improving utilization, and this improvement is duly achieved by SFSJ.

6.1.2 | Job flow

How does node stealing impact flows on the platform? First, on a supercomputer, job flows highly depend on the job sizes (i.e., number of requested compute nodes) and on the requested wall times. Indeed, even if the main scheduling algorithm is based on a First-Come-First-Served strategy: (i) backfilling strategies allow to schedule faster “small” jobs that can fit in holes of the schedule; (ii) large jobs are more frequently subject to failures. In Figure 4, we plot the response time as a function of the number of requested nodes for Baseline, without and with failures for March 2018 (Mira). Similar results are shown for June 2017 (Mira) in Figure 1 and for June 2013 (Intrepid) in Figure 5. In the figure, we report the flow of Baseline without failure (red dot) and with failures (boxplot). These flows are presented as a function of the node count of the jobs ($x = [2^n, 2^m)$ means that this is the flow of jobs whose number of nodes is in the interval $[2^n, 2^m)$), and also globally (“all”, “weighted” on the right of the x axis).

In the failure-free scenario (red dots), we see the impact of backfilling on the flow of jobs: jobs with less than

TABLE 5 Statistics for the June 2017 and March 2018 Mira workloads.

(a) Percentage of time at least one node is available right after a failure

Heuristic	(Mira) 06/2017	(Mira) 03/2018	(Intrepid) 06/2013
Baseline	92.09%	93.36%	93.83%
SFSJ	90.67%	92.48%	93.70%

(b) Number of times SFSJ interrupts a job or enrolls a free node available right after a failure

	node stealing	empty node	total failures
Trace			
(Mira) June 2017	63.4	404.4	467.8
(Mira) March 2018	46.6	416.2	462.8
(Intrepid) June 2013	45.8	410.8	456.6

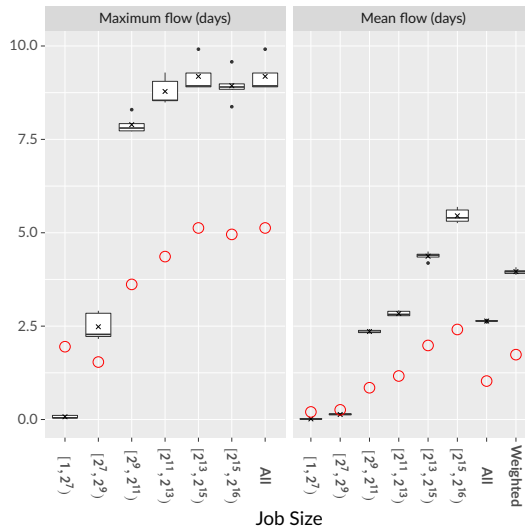


FIGURE 4 Maximum flow and mean flow as a function of job size, without failures and with failures, using Baseline. Results are for the Mira platform in March 2018.

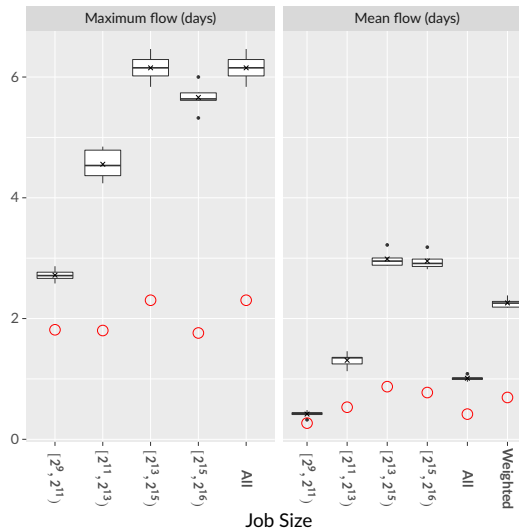


FIGURE 5 Maximum flow and mean flow as a function of job size, without failures and with failures, using Baseline. Results are for the INTREPID platform in June 2013.

$2^7 = 128$ nodes typically have a much lower flow than larger jobs. The negative impact of failures on the flows is shown by studying the difference between the failure free scenario, and the one with failures: the relative difference is much more important for larger jobs. Interestingly, failures improve the maximum flow of small jobs (jobs with less than 128 nodes). The explanation for this unexpected behavior is that failures create more "gaps" in the schedule to backfill small jobs.

With this in mind, we now compare the various flows between Baseline and SFSJ. Figures 6 and 7- show the ratio Baseline over SFSJ of the flows, hence the lower the better for SFSJ. For clarity, we also report the absolute values of these flows in Figure 8.

In these figures, we see that SFSJ significantly improves the maximum flow of large jobs, up to 10% in some scenarios, at the cost of a slight overhead in the flows of small size jobs. In the worst case, the maximum flow of small jobs is increased by a factor 2 (March 2018, Mira), but this needs to be put in perspective: the maximum flow of small jobs is several orders of magnitude lower than the flow of larger jobs. We observe that even when there is no job with a very small number of processors, as with the Intrepid-June 2013 trace, SFSJ is able to decrease the maximum flow of very large jobs, at the cost of a slight increase (larger than 20%) of the flow of medium size jobs, which is anyway at least twice smaller than that of large jobs. A similar observation is that SFSJ may slightly increase the global mean flow of the platform. Again, this is because this mean flow does not take the size of the jobs into account: if we consider the weighted mean flow instead, where the importance of the job flows depends on the node count, we do observe a decrease when using SFSJ.

Overall, SFSJ significantly improves the maximum flow of large jobs at the detriment of smaller jobs. We argue that this is a good thing since their respective absolute differ by several orders of magnitude.

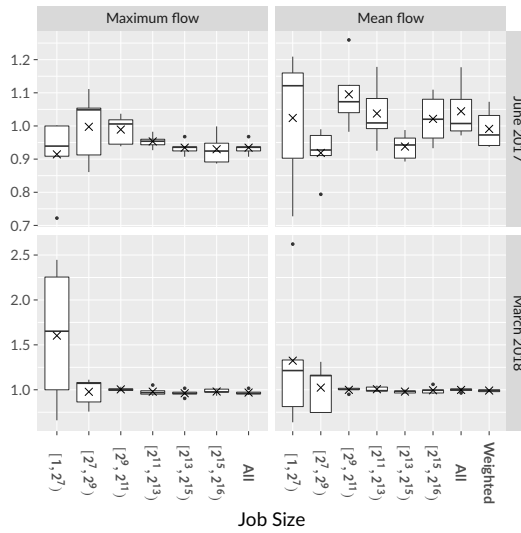


FIGURE 6 Maximum flow and mean flow ratios (SFSJ over Baseline) as a function of job sizes for Mira. For better visualization, an outlier with $x \in [1, 2^7)$, $y = 7.99$ for maximum flow, March 2018, has been hidden.

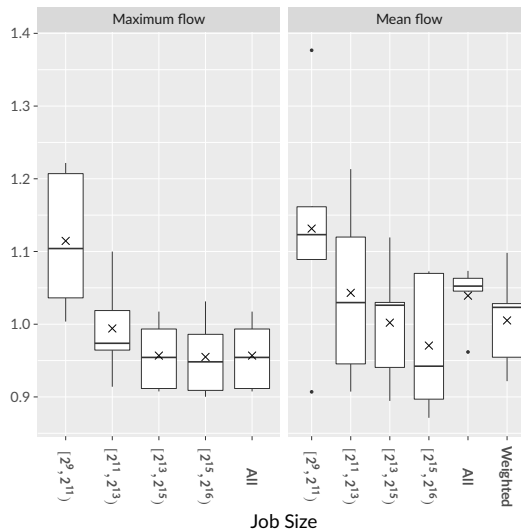


FIGURE 7 Maximum flow and mean flow ratios (SFSJ over Baseline) as a function of job sizes for Intrepid.

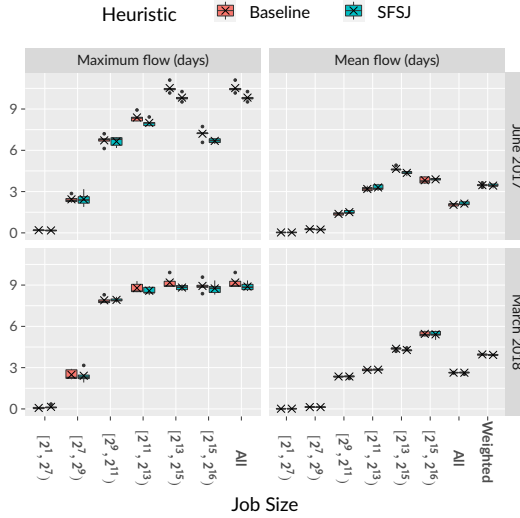


FIGURE 8 Maximum flow and mean flow as a function of job size with failures, using Baseline and SFSJ.

6.2 | Quantitative evaluation when MTBF and downtime vary

In the previous section, we have shown the positive impact of SFSJ on a given scenario. We now vary the key parameters, namely the platform MTBF and the duration of the downtime, to fully assess the usefulness of SFSJ and present its limits. We conduct experiments with MTBF $\mu \in \{20\text{min}, 40\text{min}, 1\text{h}, 2\text{h}, 5\text{h}, 10\text{h}\}$, and downtime $D \in \{10\text{min}, 1\text{h}, 1\text{day}\}$.

6.2.1 | Utilization

Figure 9 reports the ratio of the utilization of SFSJ over that of Baseline as a function of the MTBF, and for several downtime values. A value of 1.05 means that SFSJ improves the utilization by 5%, a value of 0.95 that it decreases it by 5%. From a MTBF perspective, the smaller the MTBF (i.e. the more frequent the failures are), the higher the utilization of SFSJ. Similarly, the smaller the downtime, the higher its gain in utilization. With a brief downtime (10min), the improvement of SFSJ is between 2% and 4%, while with a large downtime (1 day), its gain is negligible. This is extremely promising for future supercomputers, whose MTBF decreases linearly with size but whose downtime can (hopefully) be kept at low values.

There is one scenario where we observe a 1% loss in the utilization of SFSJ: June 2017, MTBF of 10h, downtime of 10min. This is a scenario where there are extremely few failures. When there is one, its impact is extremely small (downtime of 10 min) compared to the order of magnitude of the restart time (5 min). Hence in this scenario, killing a small job which does not expect to be killed hurts the system.

To conclude, the more failures, and the smaller the downtime, the more positive impact SFSJ has on platform utilization of the machine. There are some limit scenarios where it may be detrimental (essentially when there are few failures with a small downtime).

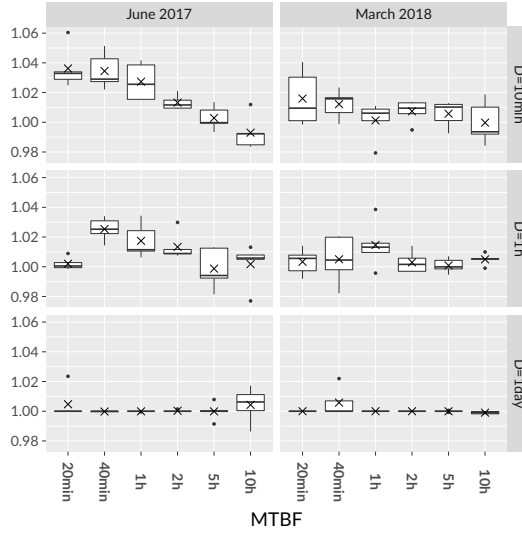


FIGURE 9 Normalized useful utilisation as a function of MTBF with failures. The higher, the better the performance of SFSJ.

TABLE 6 There are n_i jobs with $p_i = 2^i$ processors for $i \in \{0, 6\}$ in the synthetic workload.

Number of jobs n_i	504	198	108	65	55	42	28
Job size p_i	1	2	4	8	16	32	64

6.2.2 | Job flow

In Figure 10(a), we report the maximum flow of the largest jobs (jobs with more than 2^{15} nodes), and the weighted average flow (overall jobs) as a function of the MTBF for the different scenarios. In Figure 10(b), we report the ratio of these values (flow of SFSJ over that of Baseline), hence a value of 1.05 means in this figure indicates that SFSJ increases the flow by 5%, while a value of 0.95 decreases it by 5%. Both figures are for June 2017. The corresponding results for March 2018 are reported in Figure 11.

SFSJ has positive impact on both the maximum flow of the largest jobs and the weighted mean flow over all scenarios. With a small downtime (10 minutes) and a MTBF lower than 5h, the maximum flow improves up to 10-20%. This improvement is not so consequent when the downtime increases, and close to zero when the downtime is equal to 1 day.

6.3 | Evaluation with synthetic workload

In this final section we generate a synthetic workload which allows to create a different job mix. We consider a platform of 128 nodes. We create a set of 1000 jobs with various node requirements between 1 and 64 (see Table 6 for details). For application $A_j, j \in \{1, \dots, 1000\}$, we compute its delay t_{exec}^j uniformly at random between 1 minute and 119 minutes (hence an expected time of 60 minutes). Its walltime is $t_{walltime}^j = \alpha_j t_{exec}^j$, where α_j is selected uniformly at random between 1 and 5.

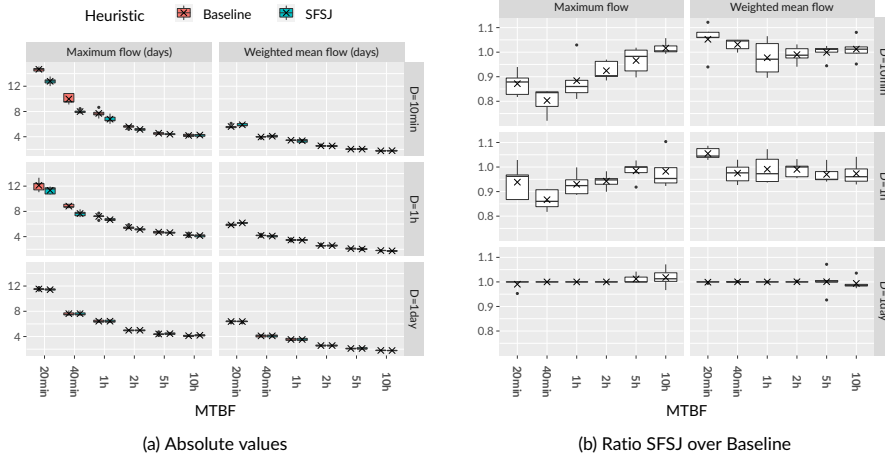


FIGURE 10 Maximum flow for largest jobs and weighted mean flow for all jobs when the MTBF and downtime vary, for June 2017. Absolute values are reported in the top plot while ratios (SFSJ over Baseline) are reported in the bottom plot.

TABLE 7 Utilization for synthetic trace.

Heuristic	Synthetic
Baseline	70%
SFSJ	72%

Jobs arrive in the system in a random order, with an inter-arrival time following an exponential distribution of mean λ (Poisson process). The value of λ is set to $174s$ so that the stress of the system is equal to 95%, i.e.:

$$\frac{1}{\lambda} \cdot \frac{\sum_i n_i \cdot 2^i \cdot 3600}{128 \cdot 1000} = 0.95. \quad (3)$$

Finally, we consider a platform MTBF of 30 minutes, a checkpoint and recovery time of 5 minutes and a downtime of 10 minutes.

Results are presented in Table 7 and Figure 13. They are similar to those observed in the previous sections: we observe an improvement of 10-15% on the maxflow and meanflow of large jobs. As a trade-off, the maximum flow of the smallest jobs increases by a large factor. However, the flow of small jobs in absolute value is much lower than that of the largest jobs (see Figure 12).

To conclude this experimental evaluation, SFSJ has positive impact on both the maximum flow of large jobs and the platform utilization of the machine, as soon as failures are not too infrequent (the very framework for which SFSJ is introduced). The impact is greater when the downtime is small.

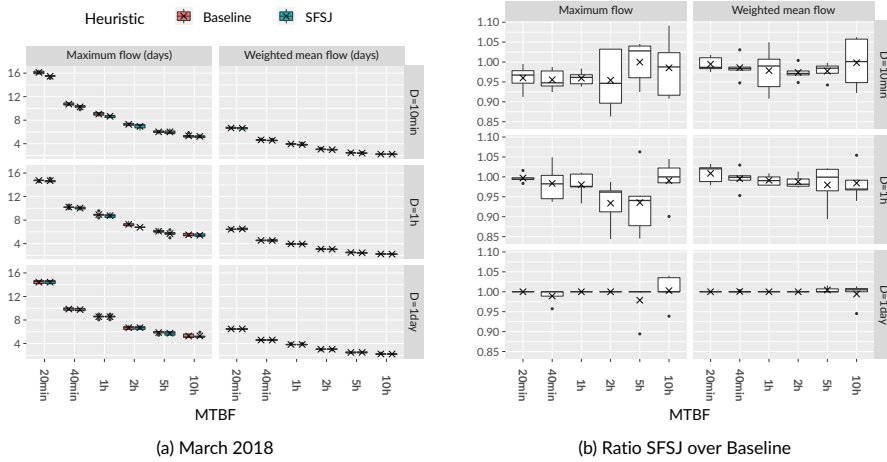


FIGURE 11 Maximum flow for largest jobs and weighted mean flow for all jobs when the MTBF and downtime vary, for March 2018. Absolute values are reported in the top plot while ratios (SFSJ over Baseline) are reported in the bottom plot.

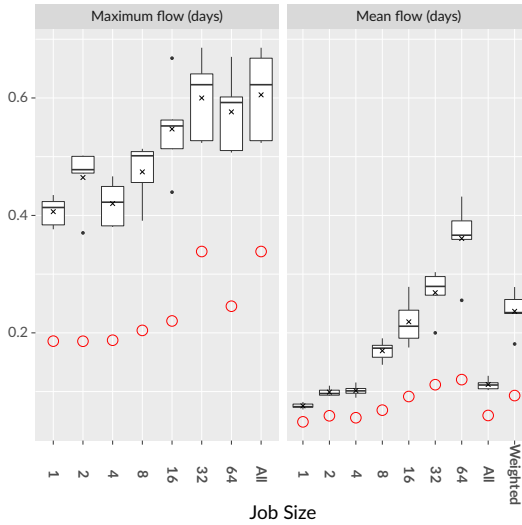


FIGURE 12 Maximum flow and mean flow for Baseline for the synthetic workload as a function of job size, without failures and with failures.

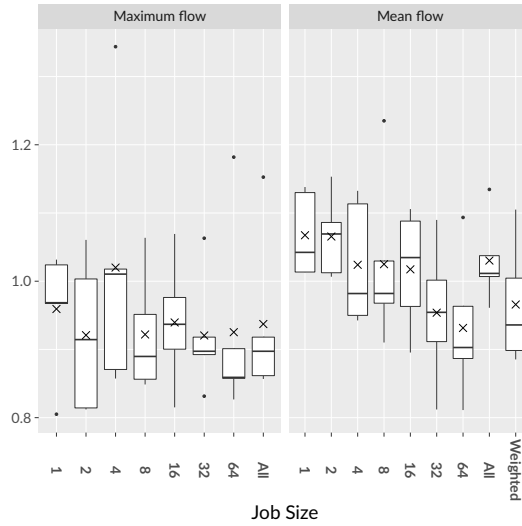


FIGURE 13 Maximum flow and mean flow ratios (SFSJ over Baseline) as a function of job sizes for synthetic workload.

7 | ADDITIONAL HEURISTICS

In Section 6, we have focused on SFSJ, the node stealing heuristic which interrupts the job with the fewest nodes, and which has been running for the smallest amount of time if there is a tie. We have designed and implemented many other variants, and compared their performance with SFSJ. These variants are sketched in Section 7.1. Details on their implementation are available in Section 7.2. Finally, the results of the simulations are reported in Section 7.3.

7.1 | Design of node-stealing variants

The first question to deal with when studying node stealing is the choice of the victim job J_{victim} , that is, which job should be considered to be interrupted to free a node so that the failed job J_{failed} can be restarted. We consider here two possible choices:

- V1** Select the currently running job with the smallest number of nodes as the victim job J_{victim} . If ties, choose the one whose release date is the latest (this is the solution chosen in Section 4.2 and evaluated in Section 5). The intuition for stopping small jobs is that they already have the smallest flows, and they are easy to reschedule.
- V2** Select the currently running job with the latest release time as the victim job J_{victim} . If ties, choose the one whose number of nodes is smallest. The idea here is to stop jobs whose waiting times are among the smallest.

Once a victim is chosen, we need to decide when we will interrupt the victim job. We propose three scenarios for this timing decision:

- T1** We immediately interrupt the victim job, and restart it from its previous checkpoint (this is the solution chosen in Section 4.2 and evaluated in Section 5).

- T2** We proactively start a checkpoint on job J_{victim} , and stop this job right after the checkpoint. This avoids wasting computation time on J_{victim} , but induces some delay for the failed job J_{failed} as it can only be restarted after the checkpoint of J_{victim} .
- T3** We wait for the next regular checkpoint of J_{victim} , and stop this job right after the checkpoint. This has a minimal impact on J_{victim} but induces a large waiting time for J_{failed} .

Finally, we need a criterion to decide it is worth interrupting it, or if we should rather wait for job to terminate to get a new node. We propose three criterion for this decision:

- K1** If the victim job J_{victim} uses strictly less nodes than the failed job J_{failed} , we decide to interrupt J_{victim} (this is the decision taken in Section 4.2).
- K2** If the victim job J_{victim} was released more recently than failed job J_{failed} , we interrupt J_{victim} .
- K3** We compute the flow of both jobs J_{victim} and J_{failed} based on their walltime in both scenarios (interrupting J_{victim} or waiting for a job completion), and we select the scenario that leads to the smallest maximum flow for these two jobs.

On the whole, we thus get 18 variants of node stealing, which are denoted by $Hxyz$, where x corresponds to timing choice T_x , y corresponds to victim choice V_y and z corresponds to interrupting choice K_z . For example SFSJ, the node stealing heuristic studied in the main paper, is denoted by H111.

7.2 | Details on the implementation

In this section, we provide details on the implementation of the heuristics. We first give some insights on the implementation of the simulations (Section 7.2.1), then we detail how to compute the remaining part of the victim job in the case of timing decision T2 and T3 (Section 7.2.2). We assume here that no new failure occurs until we have completely handled the current one, that is, until a checkpoint is taken and the failed job can be restarted. We finally explain how to handle the infrequent events of consecutive failures in Section 7.2.3.

7.2.1 | Simulation details

Algorithm 1 gives a precise statement of the various node stealing variants. Note that whenever a job is struck by a failure, its re-execution is submitted with priority 3. When a job is selected as a victim and interrupted, its re-execution is submitted with priority 2. Regular jobs have priority 1. In case of a failure (with or without node stealing), or in case of the rejuvenation of a node, the whole schedule is cleared and all jobs are rescheduled, by decreasing priority.

For timing variants T2 and T3 (proactive checkpointing and next checkpoint), the failed job is not resubmitted immediately after its failure. To avoid small jobs taking advantage of the nodes left idle by the failed jobs (that will be used for its re-execution), we submit a fictitious job to wait for the termination of the checkpoint on the victim job. If the failed job originally enrolled n nodes and had a single failure, this fictitious job uses $n - 1$ nodes. In case the failed job has no more remaining nodes (after one or multiple failures), then we can not submit this fictitious job. Because this fictitious job is needed to trigger the end of the proactive/future checkpoint on the victim job in our simulations, we cannot use node stealing in this situation. We thus cancel node stealing for this failed job and simply resubmit it from its last checkpoint. However, note that this concerns a very limited number of real scenarios.

Algorithm 1 General framework for heuristic $Hxyz$. Note that regular jobs have priority 1.

Require: Failed job J_{failed} , failure time t

```

if there is an available node  $P$  at time  $t$  then                                 $\triangleright$  No need for node stealing
    Submit the remaining part of  $J_{failed}$  (from its previous checkpoint) with priority 3 and allocate it to the previously
    allocated nodes that are still available and  $P$ 
else
    Choose the victim  $J_{victim}$  according to victim choice  $Vy$ 
    if No victim is found or the interrupting criterion  $Kz$  is negative then         $\triangleright$  Do not use node stealing
        Submit the remaining part of  $J_{failed}$  (from its previous checkpoint) with priority 3
    else                                                                             $\triangleright$  Use node stealing
        switch Timing choice  $Tx$  do
            case  $T1$                                                                      $\triangleright$  Interrupt victim right away
                Interrupt victim job  $J_{victim}$ 
                Submit the remaining part of  $J_{failed}$  with priority 3
                Submit the remaining part of  $J_{victim}$  with priority 2
            case  $T2$                                                                      $\triangleright$  Proactively checkpoint victim
                Suspend job  $J_{victim}$  and initiate a checkpoint and wait for its completion
                Interrupt  $J_{victim}$ 
                Submit the remaining part of  $J_{failed}$  with priority 3
                Submit the remaining part of  $J_{victim}$  with priority 2
            case  $T3$                                                                      $\triangleright$  Wait for next regular checkpoint
                Wait for the completion of the next checkpoint of  $J_{victim}$ 
                Interrupt  $J_{victim}$ 
                Submit the remaining part of  $J_{failed}$  with priority 3
                Submit the remaining part of  $J_{victim}$  with priority 2

```

7.2.2 | Computing the remaining part of the victim job

In all heuristics, the failed job is restarted from its previous checkpoint. The computation of the remaining part of the failed job has already been presented in section 4.1. The job selected as the victim of the node stealing either (i) can be interrupted right away, as defined by timing decision T1 (in this case, the same formulas are used to compute the characteristics of its resubmission), or (ii) it can be interrupted later for timing decisions T2 and T3: we either proactively trigger a checkpoint, or wait for the next regular checkpoint. This requires to change the computation of the characteristics of the resubmitted victim job. We now details these computations.

Proactive checkpoint (timing decision T2)

We consider here a victim job with a checkpoint period T , a checkpoint time C . In the proactive checkpoint, we may interrupt a job in the middle of a regular period, for example after a time $T_1 < T$ after the beginning of the period. Hence, when restarting the job, the first period may be different from the following ones, as it consists in completing the remaining part of this period, of length $T - T_1$. To deal with such cases, we denote by $T_{first} = T - T_1$ the duration of the first period. Since the job may be a resubmission of a previously failed or stopped job, we denote by R_{first} its initial recovery time. We have two cases:

- $R_{first} = 0$ in case of an initial submission,
- $R_{first} = R$ in case of a resubmission.

We consider that the victim job has a length t_{exec} and was started at t_{start} . The failure (on the failed job) happens at time t_{fail} . To simplify, we denote that $t_{run} = t_{fail} - t_{start}$ the length of the victim job up to the failure and $t_{first} = R_{first} + T_{first} + C$ the length of its first period, as it may differ from the following ones if the victim job is itself a resubmission of previously interrupted job. If the victim job is an initial submission, we just let $T_{first} = T$ and $R_{first} = 0$. In Figure 14, we illustrate these notations on the execution of a job. We will use times t_1, \dots, t_5 as potential times for failures in the description of the formulas below.

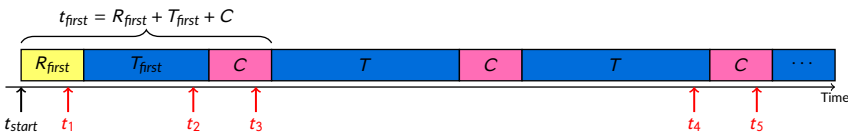


FIGURE 14 Illustration of the notations for the victim job, with the five cases distinguished to compute the proactive checkpoint strategy (T2).

We aim at computing when the victim job is interrupted and what are the characteristics of its resubmission. More precisely, we will first compute the following two quantities:

- The time t_{useful} spent by the victim job doing useful work until we stop it. This includes execution time and regular checkpoint time, but not the checkpoint that we introduce due to the proactive checkpoint strategy.
- The time $t_{checkpoint}$ need to complete the checkpoint introduced by the proactive checkpoint strategy. This checkpoint will be completed at time $t_{fail} + t_{checkpoint}$: at this time we will be able to steal a node from the victim to restart the failed job.

Then we will compute the characteristics of the resubmitted victim job: its length t'_{exec} , its first period length T'_{first} and its recovery time R'_{first} .

We distinguish five cases depending on which part of the job is struck by a failure.

The first case is when the failure occurs during the execution of the potential recovery R_{first} . This is when the failure hits at time $t_{fail} = t_1$ in the Figure 14. This means that no progress was made in the job ($t_{useful} = 0$), hence there is no need to start a proactive checkpoint and we can simply resubmit the job without any modifications. Note that this case only happens if the victim job was a re-execution of a job (as $R_{first} > 0$), so we have $R_{first} = R$, and $R'_{first} = R$,

The second case is when the failure occurs during the execution of the first period T_{first} , which corresponds to $t_{fail} = t_2$ in Figure 14. Then we start a proactive checkpoint at time t_{fail} to save the useful work executed in T_{first} . In this case, the saved useful work is $t_{useful} = t_{run} - R_{first}$ (since we classify regular checkpoint into the useful work, and proactive checkpoint not into the useful work). In the resubmission job, the first period will have to terminate the execution of the interrupted first period, that is, it will run from t_2 to T_{first} . Hence we will set $T'_{first} = T_{first} - (t_{run} - R_{first})$. The time needed to complete the proactive checkpoint is $t_{checkpoint} = C$.

The third case is when the failure occurs during the checkpoint that follows the first period, as illustrated by $t_{fail} = t_3$ in Figure 14. Then we do not need to start a proactive checkpoint, we simply wait for the completion of the ongoing checkpoint to be completed. The time we have to wait for the completion of the checkpoint is $t_{checkpoint} = t_{start} + t_{first} - t_{fail} = t_{first} - t_{run}$. In this case, the useful work performed by the job is $t_{useful} = T_{first} + C$ (containing the regular checkpoint into the useful work). The resubmission of the victim job will start by a regular period, that is, $T'_{first} = T$.

The fourth case happens when the failure occurs during a regular period T , which corresponds to $t_{fail} = t_4$ in Figure 14. We then start a proactive checkpoint at time t_{fail} to save the work already performed in this period, hence the duration of this checkpoint is $t_{checkpoint} = C$. The amount of successful work is thus $t_{useful} = t_{run} - R_{first}$. The first period of the resubmitted job copy will perform the missing work from t_{fail} to the end of the regular period, computed as

$$T'_{first} = T - (t_{run} - t_{first} - \left\lfloor \frac{t_{run} - t_{first}}{T + C} \right\rfloor \times (T + C)).$$

The fifth case happens when the failure occurs during a regular checkpoint, for example for $t_{fail} = t_5$ in Figure 14. As in the third case, we do not start a proactive checkpoint but take advantage of the ongoing one. In this case, the useful work starts at the beginning of T_{first} until the end of this regular checkpoint, that is

$$t_{useful} = T_{first} + C + \left\lfloor \frac{t_{run} - t_{first}}{T + C} \right\rfloor \times (T + C).$$

The time we have to wait until the end of the current checkpoint goes from t_{fail} until the end of the ongoing checkpoint, that is

$$t_{checkpoint} = T + C - (t_{run} - t_{first} - \left\lfloor \frac{t_{run} - t_{first}}{T + C} \right\rfloor \times (T + C)).$$

In this case, the first period of the resubmitted victim job is a regular one, so that $T'_{first} = T$.

In cases 2 to 5, we start proactive checkpoints, so that the recovery time for the resubmission of the victim job is set to $R'_{first} = R$.

Figure 15 summarizes how to compute the length of the resubmission of the victim job, as well as the time

$$\begin{aligned}
t_{useful} &= \begin{cases} 0, & \text{if } t_{run} \leq R_{first}, \\ t_{run} - R_{first}, & \text{if } R_{first} < t_{run} \leq R_{first} + T_{first}, \\ T_{first} + C, & \text{if } R_{first} + T_{first} < t_{run} \leq t_{first}, \\ t_{run} - R_{first}, & \text{if } t_{run} - t_{first} - \left\lfloor \frac{t_{run} - t_{first}}{T+C} \right\rfloor \times (T+C) \leq T, \\ T_{first} + C + \left\lfloor \frac{t_{run} - t_{first}}{T+C} \right\rfloor \times (T+C), & \text{otherwise.} \end{cases} \\
T'_{first} &= \begin{cases} T_{first}, & \text{if } t_{run} \leq R_{first}, \\ T_{first} - (t_{run} - R_{first}), & \text{if } R_{first} < t_{run} \leq R_{first} + T_{first}, \\ T, & \text{if } R_{first} + T_{first} < t_{run} \leq t_{first}, \\ T - (t_{run} - t_{first} - \left\lfloor \frac{t_{run} - t_{first}}{T+C} \right\rfloor \times (T+C)), & \text{if } t_{run} - t_{first} - \left\lfloor \frac{t_{run} - t_{first}}{T+C} \right\rfloor \times (T+C) \leq T, \\ T, & \text{otherwise.} \end{cases} \\
t_{checkpoint} &= \begin{cases} 0, & \text{if } t_{run} \leq R_{first}, \\ C, & \text{if } R_{first} < t_{run} \leq R_{first} + T_{first}, \\ t_{first} - t_{run}, & \text{if } R_{first} + T_{first} < t_{run} \leq t_{first}, \\ C, & \text{if } t_{run} - t_{first} - \left\lfloor \frac{t_{run} - t_{first}}{T+C} \right\rfloor \times (T+C) \leq T, \\ T + C - (t_{run} - t_{first} - \left\lfloor \frac{t_{run} - t_{first}}{T+C} \right\rfloor \times (T+C)), & \text{otherwise.} \end{cases} \\
R'_{first} = R, \quad t'_{exec} = R'_{first} + t_{exec} - t_{useful}, \quad t'_{walltime} = R'_{first} + t_{walltime} - t_{useful}. & \quad (4)
\end{aligned}$$

FIGURE 15 Formulas used to compute the useful executed work of the victim job t_{useful} , the length of the first period of the resubmitted victim job T'_{first} , the time between the failure and the completion of the proactive checkpoint $t_{checkpoint}$, the recovery time of the resubmitted victim job R'_{first} , and the length t'_{exec} and wall time $t'_{walltime}$ of the resubmitted victim job.

$t_{checkpoint}$ to wait until a node can be given to the failed job to restart it.

Using next regular checkpoint (timing decision T3)



FIGURE 16 Illustration of the notations for the victim job, with the five cases distinguished to compute the future checkpoint strategy (T3).

We now detail how to compute the remaining time of the victim job as well as the time when a new node is available for the failed job in the case of the future checkpoint heuristic (timing decision T3). We use the definition introduced above for the proactive checkpoint heuristic. In the future checkpoint heuristic, all periods between checkpoints have the same length, contrarily to the proactive checkpoint heuristic when the first period may be different from the other. This means we always have $T_{first} = T$. This largely simplifies the analysis. We now distinguish between two cases depending on the state of the victim when the failure happens, illustrated on Figure 16.

The first case happens when the failure occurs during the execution of R_{first} by the victim job. This is the case for example for $t_{fail} = t_1$ in Figure 16. Then, there is no reason to wait for the next regular checkpoint, as no useful work has been performed by the victim job. We simply interrupt the victim and submit it again later. We thus have $t_{useful} = 0$ and $t_{checkpoint} = 0$.

The second case happens when the failure occurs after R_{first} , either during a regular execution or during a checkpoint, as illustrated by $t_{fail} = t_2$ or $t_{fail} = t_3$ on Figure 16. Then we need to wait until the next regular checkpoint of the victim job. In this case, the time performing useful work starts at the end of R_{first} and goes to the completion of the next checkpoint, that is, the one immediately following t_{fail} (we recall that regular checkpoints are counted as useful work). Hence we have

$$t_{useful} = \left\lceil \frac{t_{run} - R_{first}}{T + C} \right\rceil \times (T + C).$$

The time between the failure at t_{fail} and the completion of the next checkpoint can be computed as:

$$t_{checkpoint} = R_{first} + \left\lceil \frac{t_{run} - R_{first}}{T + C} \right\rceil \times (T + C) - t_{run}.$$

Again, the first case only happens when $R_{first} > 0$, that is $R_{first} = R$. Hence, in both cases, the recovery time of the victim job copy is $R'_{first} = R$.

To sum up, we compute the useful working time for the victim job as follows:

$$t_{useful} = \begin{cases} 0, & \text{if } t_{run} \leq R_{first}, \\ \left\lceil \frac{t_{run} - R_{first}}{T + C} \right\rceil \times (T + C), & \text{otherwise.} \end{cases}$$

The time that we need to wait between t_{fail} and the completion of the next checkpoint of the victim job (which

is the delay of the failed node needs to wait before being restarted with a stolen node) is computed as follows

$$t_{checkpoint} = \begin{cases} 0, & \text{if } t_{run} \leq R_{first}, \\ R_{first} + \left\lceil \frac{t_{run} - R_{first}}{T+C} \right\rceil \times (T+C) - t_{run}, & \text{otherwise.} \end{cases}$$

The length of the resubmission of the victim job is finally be computed as previously, as well as its wall time:

$$\begin{aligned} t'_{exec} &= R'_{first} + t_{exec} - t_{useful}, \\ t'_{walltime} &= R'_{first} + t_{walltime} - t_{useful}. \end{aligned}$$

7.2.3 | Consecutive failures

In the previous discussion, we assumed that no failure hits either the victim job or the remaining node of the failed job until the checkpoint is completed and the failed job may be restarted. However, such rare cases can happen. We detail here how to handle them.

We assume that a job failed because of a failure at time t_{fail}^{first} . A victim job was selected in order to perform node stealing, that is, to relaunch the failed job with its remaining nodes plus one node of the victim job. In timing decision T2, we trigger a proactive checkpoint on the victim job at time t_{fail}^{first} as shown in Figure 17. In timing decision T3, we wait until the next regular checkpoint of the victim job, as shown in Figure 18. We now consider the event of a failure before the end of the checkpoint, on the victim job or on the nodes of the failed job that were not hit by a failure at time t_{fail}^{first} and remained available.

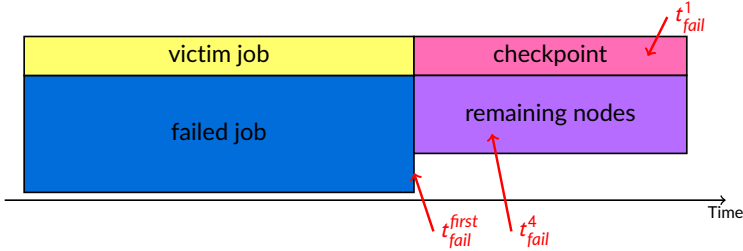


FIGURE 17 Notations of the special case for heuristic T2.

The first case is that old victim job or its proactive or regular checkpoint is hit by a new failure, as shown by t_{fail}^1 , t_{fail}^2 and t_{fail}^3 in Figures 17 and 18. When this happens, we cancel node stealing for the originally failed job: both jobs are simply restarted from their previously successful checkpoints.

The second case deals with a new failure striking the remaining nodes of the originally failed job, as shown by t_{fail}^4 and t_{fail}^5 in Figures 17 and 18. In this case, there are two possibilities:

- The number of nodes of the victim job is smaller than the total number of failures that hit nodes of the failed job, which means the number of nodes of the old victim job is not sufficient to resume the old failed job. In this case, we also cancel node stealing and resubmit the failed job from its previous successful checkpoint. Since node stealing is canceled, the victim job continues its execution until its regular termination.

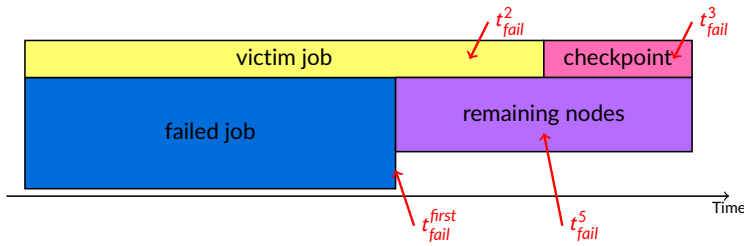


FIGURE 18 Notations of the special case for heuristic T3.

- The number of nodes of the victim job is larger than or equal to the number of failures occurred on nodes of the failed job, which means the number of nodes of the victim job is enough to resume the old failed job. In this case, we continue using node stealing for the failed job with the same victim job.

Another special case may occur when using future checkpoint (timing decision T3): the victim job may complete before its next checkpoint. In this case, we simply resubmit the failed job right after the termination of the victim job.

7.3 | Presentation of all results and discussion

In this section, we report the results of experiments comparing the different variants of node stealing introduced above.

Utilization

We start by comparing the useful utilization of the platform by all heuristics, as presented in Table 8 which generalizes Table 4. This table also present the percentage of time at least one spare node is available (as previously in Table 5). In this table (and below), we recall that variant $x\ yz$ denotes the algorithm obtained with timing choice Tx, victim choice Vy and interrupting choice Kz.

We first remark that no variant is able to really increase the utilization above what is achieved by the initial node stealing heuristic (variant 111). Some of them even decrease the utilization below the one achieved by the baseline heuristic by up to 5%. As previously, we relate the impact on utilization to the percentage of time an idle node is available as a spare (and thus, node stealing is not useful). The table clearly shows that the larger this percentage of time, the smaller the impact of node stealing.

We also measure the number of time node stealing is used in Table 9. We remark that only variant 311, (corresponding to waiting the completion of the next regular checkpoint of the victim to interrupt it) is able to increase the usage of node stealing compared to our initial proposal. Using proactive checkpointing (variants $x\ **$) can (slightly) increase or decrease the use of node stealing. The other possibilities for y (choice of the victim) and z (interrupting criterion) always lead to a reduced usage of node stealing.

Job flows

Note that in the following figures (Figure 19, 20 and 21), a few outliers has been omitted for better readability. There are described in Table 10.

We first study the effect of changing the timing decision on the performance of node stealing. Figure 19 depicts

Variant	June 2017		March 2018	
	utilization	idle perc.	utilization	idle perc.
baseline	79.32 %	92.09 %	77.39 %	93.36 %
111	80.69 %	90.67 %	78.50 %	92.48 %
112	80.15 %	91.95 %	78.12 %	92.00 %
113	80.49 %	91.85 %	77.64 %	92.65 %
121	80.48 %	91.67 %	77.47 %	92.61 %
122	79.40 %	91.50 %	76.95 %	92.63 %
123	79.20 %	91.76 %	77.15 %	92.77 %
211	80.52 %	90.87 %	78.09 %	92.01 %
212	80.30 %	90.81 %	78.11 %	92.53 %
213	80.50 %	90.35 %	77.91 %	92.18 %
221	80.50 %	91.16 %	78.14 %	92.45 %
222	79.60 %	91.22 %	77.94 %	92.46 %
223	80.17 %	91.06 %	78.31 %	92.42 %
311	76.18 %	85.87 %	75.17 %	88.90 %
312	77.50 %	89.48 %	76.41 %	90.57 %
313	78.35 %	90.32 %	76.87 %	91.38 %
321	78.82 %	89.87 %	76.88 %	92.05 %
322	78.57 %	89.03 %	75.92 %	90.09 %
323	78.57 %	89.68 %	77.28 %	92.15 %

TABLE 8 Utilization and percentage of time at least one node is available (idle perc.) in all variants. Results are for the Mira platform in June 2017 and March 2018.

Variant	June 2017			March 2018		
	node	empty	total	node	empty	total
	stealing	node	failures	stealing	node	failures
111	63.4	404.4	467.8	46.6	416.2	462.8
112	33.6	433.4	467.0	44.6	415.6	460.2
113	48.0	423.6	471.6	44.6	420.0	464.6
121	32.4	440.4	472.8	29.8	432.6	462.4
122	46.8	422.4	469.2	49.8	414.4	464.2
123	48.0	418.2	466.2	48.2	414.2	462.4
211	60.8	413.2	474.0	53.2	411.2	464.4
212	34.6	431.6	466.2	43.8	417.2	461.0
213	47.6	424.2	471.8	41.0	419.8	460.8
221	33.2	432.2	465.4	27.4	432.4	459.8
222	49.8	410.4	460.2	48.0	413.2	461.2
223	51.4	417.6	469.0	49.4	411.8	461.2
311	76.4	362.6	439.0	60.2	383.6	443.8
312	39.8	413.4	453.2	43.6	408.4	452.0
313	17.2	437.4	454.6	14.8	441.2	456.0
321	32.8	427.4	460.2	26.2	430.4	456.6
322	55.8	395.8	451.6	57.0	394.0	451.0
323	30.6	429.4	460.0	24.0	435.4	459.4

TABLE 9 Number of time node stealing is used vs number of time an empty node is used for all node stealing variants. Results are for the Mira platform in June 2017 and March 2018.

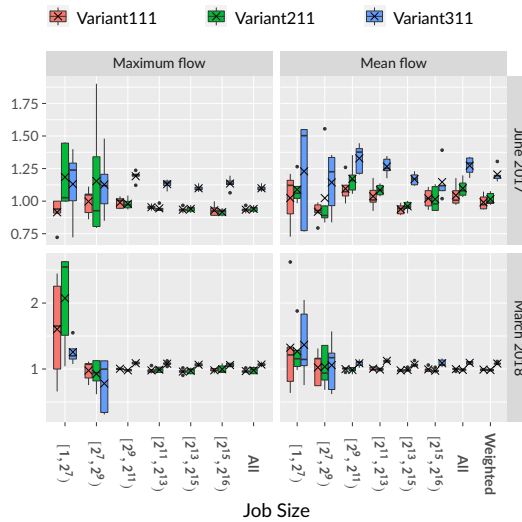


FIGURE 19 Maximum flow and mean flow relative to Baseline for various timing decision (111: immediately interrupting, 211: proactive checkpointing, 311: waiting next checkpoint) and for various categories of job sizes. MTBF and downtime are both set to 1 hour.

the results when triggering proactive checkpoint, or when using the next regular checkpoint, rather than interrupting the victim as soon as possible. We notice that no strategy is able to clearly outperform the original node stealing heuristic. Waiting for the next checkpoint always increases the maximum and average flows. Using proactive checkpointing has comparable performance with the original node stealing for large jobs, but sometimes largely increases the maximum flow of small jobs.

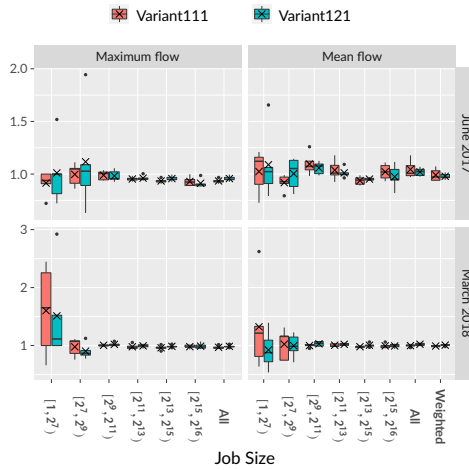


FIGURE 20 Maximum flow and mean flow relative to Baseline heuristic for the two victim choices (111: fewest nodes, 121: latest release date) and for various categories of job sizes. MTBF and downtime are both set to 1 hour.

We then study on Figure 20 the impact of the choice of the victim on the performance of node stealing. In the original node stealing, we select the job with the smallest number of nodes. In the proposed variant, we select the job with the latest release date. We see that the victim selection policy has a limited impact for small jobs, but very little for large jobs. On the whole, it does not allow to improve performance.

Finally, Figure 21 presents the results when changing the interrupting criterion. In the original heuristic, we decide to interrupt a victim job and perform node stealing if the victim requires less nodes than the failed job. We also proposed to use release date to take this decision, by interrupting a victim only if it was release later than the failed job. The last criterion requires to compute an estimation of the flow for both the failed and victim job: the victim is interrupted only if it leads to a smaller maximum flow for both jobs. We notice in these results that changing the interrupting criterion has an impact only on small jobs, and does not clearly improve the results.

On the whole, all proposed variants fail to clearly improve the performance of node stealing: the basic node stealing heuristic is sufficient to improve the flow of large jobs, at the cost of a limited increase in the flow of small jobs (which is originally much smaller than the one of large jobs) This corroborates the analysis of the utilization conducted in Section 6.1.1.

8 | RELATED WORK

In this section, we give a few pointers to related work on batch schedulers (Section 8.1) and on fault-tolerance mechanisms (Section 8.2).

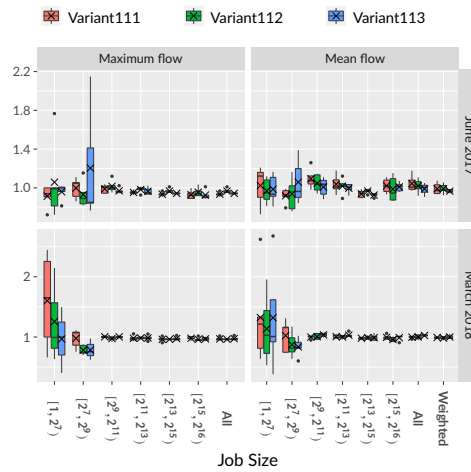


FIGURE 21 Maximum flow and mean flow relative to Baseline for the various interrupting criterion (111: less nodes, 112: later release date, 113: better estimated maximum flow) and for various categories of job sizes. MTBF and downtime are both set to 1 hour.

	Variant	job size	value	type of figure
Figure 19	111	$[1, 2^7]$	7.99	Maxflow
	311	$[1, 2^7]$	6.09	Maxflow
Figure 20	111	$[1, 2^7]$	7.99	Maxflow
Figure 21	111	$[1, 2^7]$	7.99	Maxflow
	112	$[1, 2^7]$	5.47	Maxflow
	113	$[1, 2^7]$	4.09	Maxflow
	113	$[1, 2^7]$	4.84	Maxflow
	112	$[1, 2^7]$	7.73	Maxflow
	112	$[1, 2^7]$	3.58	Meanflow

TABLE 10 Outliers removed from Figures 19, 20 and 21 for the March 2018 dataset.

8.1 | Job management on HPC platforms

Resource and Job Management Systems (RJMS), a.k.a. Batch schedulers, are intermediary software layers generally managed by a system administrator (examples include Slurm [27], Moab/Maui [16], OAR [4] etc). This software is in charge of allocating the different jobs through a scheduling heuristic, while taking into account various constraints. The most important constraints that a batch scheduler has to account for are the *estimation* by users of the resources needed for a job, both in a spatial dimension (number of processing units), and in a temporal dimension (estimated processing time).

Natural developments in batch scheduling have included more dimensions to the scheduling heuristic, such as heterogeneity of computing resources, fairness to deal with the disparity of job requirements and usage, etc. The scheduling heuristics are typically implemented by the introduction of specific *queues*, where jobs with similar characteristics (size, reservation length, priority, ...) are grouped together into the same queue [28]. Each queue is configured with a specific scheduling heuristic.

There are several main scheduling heuristics used for batch scheduling. The default for most schedulers is the First-Come-First-Served policy (FCFS) [26, 16]. This strategy is often tweaked by including the time of arrival (i.e. "first-come" condition) into a more general priority-based, greedy heuristic, that includes a wide range of parameters [26]. Other common strategies exist such as Smallest Job First [12], known to be efficient with respect to the response time objective. The advantage of these greedy heuristics is their low scheduling cost. Their drawback is a less efficient solution with a lot of idle time for the platform. To mitigate this limitation, these heuristics are coupled with a *backfilling* strategy. Backfilling consists in scheduling small jobs in the gaps created by the scheduling solutions. The two main flavors of backfilling are *conservative* (no job in the queue can be delayed by a backfilled job) and *EASY* (the first job in the queue is never delayed by backfilled jobs) [20]. This work has focused on using the conservative approach, but we expect that using EASY or other approaches would lead to very similar results and conclusions.

8.2 | Fault-tolerance from a system perspective

To mitigate the impact of node crashes, several techniques are considered, such as replication and checkpointing. In this work we consider the de-facto standard approach for HPC, periodic checkpointing [13]. With this technique, users are invited to checkpoint their jobs periodically, with the idea that if a node crashes during execution, then the job will be able to resume from the last checkpoint, instead of resuming from scratch. A key advantage of checkpointing is to decrease the amount of re-executed work after a crash. One must decide how often to checkpoint, i.e., derive the optimal checkpointing period. An optimal strategy is defined as a strategy that minimizes the expectation of the execution time of the application. For a preemptible application, i.e., an application that can be checkpointed at any time-step, the classical formula due to Young [31] and Daly [7] states that the optimal checkpointing period is $P_{\text{VD}} = \sqrt{2\mu C}$, given a checkpointing cost C and platform MTBF μ .

However, there are several complications related to deciding when, and on which resources, the job will be allowed to resume execution after experiencing the loss of one node. Several batch schedulers [15] will reschedule a failed job with high priority, thereby enabling an immediate re-execution if there is a free node available. The high priority allows the failed job to avoid a long wait in the job submission queue. Without priority, the delay between the interruption of a job and the beginning of its re-execution is called the resubmission time. Its value typically ranges from several hours to several days if the platform is over-subscribed (up to 10 days for large jobs on the K -computer [30]).

Hori et al. [14] discussed how one can use spare nodes to restart an application which has experienced a node

failure. This technique could be applied to our case. In [14], spare nodes are reserved and used only in the case of a failure, which enables the failed job to restart as fast as possible. Prabhakaran et al. [23] discussed the limitations of the reservation of spare nodes, which creates a non-negligible overhead. Instead they study the case where jobs are moldable and/or malleable; when no idle nodes is available for a failed job to restart, they propose several strategies such as executing the failed job on fewer nodes, or to take a node from a malleable job. In contrast, our work applies to rigid nodes (neither moldable nor malleable) and never changes the size of the jobs.

Recently, Fan et al. [11] have discussed the possibility of killing rigid jobs. This is one of several strategies in [11] to deal with the arrival of "on demand" jobs in the presence of hybrid workloads (rigid, moldable and malleable jobs). The arrival of these "on demand" jobs could be seen as the preemption requested by the failure of another job in our model. Interestingly, their observations (e.g., Observation 13 in [11]) seems to imply that it would be better to checkpoint more frequently rigid jobs, when our observations imply that the waste due to interruption is negligible compared to the waste due to checkpointing (and we should not checkpoint more frequently). However, checkpointing may have been included in their measure of system utilization (contrarily to this paper) which would explain the difference in observations.

The optimization of fault-tolerance techniques often considers a short downtime (also called rejuvenation time) for the failed resources, compared to the platform MTBF. This makes sense when one simply needs to reboot the machine that failed. But in the case of a defective component to be replaced, the downtime can last up to one day, because maintenance is operated at a fixed time every day, e.g. every morning for the K -computer [30]. Our experiments aimed at covering the whole range of possible values for the downtime.

Finally, outside the HPC community, several recent papers investigate fault-tolerant techniques for cloud systems and platforms [18, 1, 32, 8], but none uses node stealing techniques.

9 | CONCLUSION

Patel et al wrote in their SC'20 paper [22]: *Users are now submitting medium-sized jobs because the waits times for larger sizes tends to be longer*. Indeed, we have shown that failures dramatically increase the flow of large jobs. It is important to invent scheduling strategies that decrease the flow of large jobs on large-scale machines.

We have introduced node stealing as an efficient approach to decrease the flow of large jobs. For example, in June 2017 on Mira, the maximum flow of large jobs ([32K, 64K] nodes) goes down from 7.20 to 3.72 days, while the maximum flow of small jobs ([1, 128] nodes) increases from 0.19 to 0.54 days. We argue that the sharp decrease of the flow of large jobs is well worth the small increase of the flow of small jobs, given that large-scale platforms are primarily intended to execute large jobs. A side advantage of node stealing is a slight increase in terms of platform utilization.

We have designed several variants of node stealing and report that they behave similarly. Future work will be devoted to explore other well-established batch scheduling strategies (such as EASY) and assess the usefulness of node stealing when coupled with these strategies. A long-term objective is to design a node-stealing-aware batch scheduler: when taking scheduling decisions at submission time, the goal would be to account for the possibility of mitigate a failure by node stealing.

Currently node stealing is used very few times because there is often a node available when a failure strikes a job. Hence the failed job could restart with this node instead of requiring to steal a node from another running job. This opportunity would be drastically limited on architectures where topology matters: even if a spare node is available, it may not be possible to use it because of topology constraints. In such a framework, node-stealing would probably

lead to a dramatic decrease in response time.

Acknowledgments

We would like to thank the reviewers for their comments and suggestions, which greatly helped improve the final version of the paper. This work was supported in part by the “Adaptive multitier intelligent data manager for Exascale (ADMIRE)” project, funded by the European Union’s Horizon 2020 JTI-EuroHPC Research and Innovation Programme (grant 956748).

References

- [1] Z. Ahmad, B. Nazir, and A. Umer. A fault-tolerant workflow management system with quality-of-service-aware scheduling for scientific workflows in cloud computing. *International Journal of Communication Systems*, 34(1), 2021.
- [2] P. Brucker. *Scheduling Algorithms*. Springer, 2004.
- [3] P. Brucker and S. Knust. *Complex Scheduling*. Springer, 2012.
- [4] N. Capit, G. Da Costa, Y. Georgiou, G. Huard, C. Martin, G. Mounié, 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.*, volume 2, pages 776–783. IEEE, 2005.
- [5] H. Casanova, A. Giersch, A. Legrand, M. Quinson, and F. Suter. Versatile, scalable, and accurate simulation of distributed applications and platforms. *Journal of Parallel and Distributed Computing*, 74(10):2899–2917, 2014.
- [6] S. Coghlan, K. Kumaran, R. M. Loy, P. Messina, V. Morozov, J. C. Osborn, S. Parker, K. Riley, N. A. Romero, and T. J. Williams. Argonne applications for the IBM Blue Gene/Q, Mira. *IBM Journal of Research and Development*, 57(1/2):12–1, 2013.
- [7] J. T. Daly. A higher order estimate of the optimum checkpoint interval for restart dumps. *Future Generation Comp. Syst.*, 22(3):303–312, 2006.
- [8] T. Dong, F. Xue, H. Tang, and C. Xiao. Deep reinforcement learning for fault-tolerant workflow scheduling in cloud environment. *Applied Intelligence*, 53(9):9916–9932, 2023.
- [9] Y. Du. Implementation of the node stealing scheduler in Batsched, 2022. <https://gitlab.inria.fr/yishu/node-stealing-for-resilience/-/tree/master/tmp>.
- [10] P.-F. Dutot, M. Mercier, M. Poquet, and O. Richard. Batsim: a realistic language-independent resources and jobs management systems simulator. In *Job Scheduling Strategies for Parallel Processing*, pages 178–197. Springer, 2015.
- [11] Y. Fan, Z. Lan, P. Rich, W. Allcock, and M. E. Papka. Hybrid workload scheduling on HPC systems. In *2022 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 470–480. IEEE, 2022.
- [12] M. Harchol-Balter, B. Schroeder, N. Bansal, and M. Agrawal. Size-based scheduling to improve web performance. *ACM Transactions on Computer Systems (TOCS)*, 21(2):207–233, 2003.
- [13] T. Herault and Y. Robert, editors. *Fault-Tolerance Techniques for High-Performance Computing*, Computer Communications and Networks. Springer Verlag, 2015.
- [14] A. Hori, K. Yoshinaga, T. Herault, A. Bouteiller, G. Bosilca, and Y. Ishikawa. Sliding substitution of failed nodes. In *Proceedings of the 22nd European MPI Users’ Group Meeting, EuroMPI ’15*, New York, NY, USA, 2015. Association for Computing Machinery.

- [15] IBM Spectrum LSF Job Scheduler. Fault tolerance and automatic management host failover, 2021. <https://www.ibm.com/docs/en/spectrum-lsf/10.1.0?topic=cluster-fault-tolerance>.
- [16] D. Jackson, Q. Snell, and M. Clement. Core algorithms of the maui scheduler. In *Workshop on Job Scheduling Strategies for Parallel Processing*, pages 87–102. Springer, 2001.
- [17] D. Kumar, G. Baranwal, Z. Raza, and D. P. Vidyarthi. A survey on spot pricing in cloud computing. *Journal of Network and Systems Management*, 26(4):809–856, 2018.
- [18] Z. Li, V. Chang, H. Hu, H. Hu, C. Li, and J. Ge. Real-time and dynamic fault-tolerant scheduling for scientific workflows in clouds. *Information Sciences*, 568:13–39, 2021.
- [19] L. Lin, L. Pan, and S. Liu. Backup or not: An online cost optimal algorithm for data analysis jobs using spot instances. *IEEE Access*, 8:144945–144956, 2020.
- [20] A. W. Mu'alem and D. G. Feitelson. Utilization, predictability, workloads, and user runtime estimates in scheduling the IBM SP2 with backfilling. *IEEE transactions on parallel and distributed systems*, 12(6):529–543, 2001.
- [21] V. Munhoz, M. Castro, and O. Mendizabal. Strategies for fault-tolerant tightly-coupled HPC workloads running on low-budget spot cloud infrastructures. In *2022 IEEE 34th International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD)*, pages 263–272. IEEE, 2022.
- [22] T. Patel, Z. Liu, R. Kettimuthu, P. Rich, W. Allcock, and D. Tiwari. Job characteristics on large-scale systems: long-term analysis, quantification, and implications. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–17. IEEE, 2020.
- [23] S. Prabhakaran, M. Neumann, and F. Wolf. Efficient fault tolerance through dynamic node replacement. In *2018 18th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID)*, pages 163–172. IEEE, 2018.
- [24] S. Schlagkamp, R. Ferreira da Silva, W. Allcock, E. Deelman, and U. Schwiegelshohn. Consecutive job submission behavior at Mira supercomputer. In *Proceedings of the 25th ACM International Symposium on High-Performance Parallel and Distributed Computing*, pages 93–96, 2016.
- [25] A. Silberschatz, P. B. Galvin, and G. Gagne. *Operating System Concepts*. Wiley, 2012.
- [26] Slurm team. Slurm Multifactor Priority Plugin. https://slurm.schedmd.com/priority_multifactor.html, 2020.
- [27] Slurm team. Slurm Workload Manager. <https://slurm.schedmd.com/>, 2020.
- [28] W. Tang, N. Desai, D. Buettner, and Z. Lan. Analyzing and adjusting user runtime estimates to improve job scheduling on the blue gene/p. In *2010 IEEE International Symposium on Parallel & Distributed Processing (IPDPS)*, pages 1–11. IEEE, 2010.
- [29] Top500. Top 500 Supercomputer Sites, June 2022. <https://www.top500.org/lists/top500/2022/06/>.
- [30] K. Yamamoto, A. Uno, H. Murai, T. Tsukamoto, F. Shoji, S. Matsui, R. Sekizawa, F. Sueyasu, H. Uchiyama, M. Okamoto, N. Ohgushi, K. Takashina, D. Wakabayashi, Y. Taguchi, and M. Yokokawa. The K computer operations: Experiences and statistics. In D. Abramson, M. Lees, V. V. Krzhizhanovskaya, J. J. Dongarra, and P. M. A. Sloot, editors, *Proceedings of the International Conference on Computational Science (ICCS)*, volume 29 of *Procedia Computer Science*, pages 576–585. Elsevier, 2014.
- [31] J. W. Young. A first order approximation to the optimum checkpoint interval. *Comm. of the ACM*, 17(9):530–531, 1974.
- [32] J. Zeng. An algorithm for online distributed fault-tolerant job scheduling in grid computing. *International Journal of Web and Grid Services*, 17(4):389–407, 2021.