

A Generic Approach to Scheduling and Checkpointing Workflows

Li Han
East China Normal University, China
ENS Lyon, France
Li.Han@ens-lyon.fr

Valentin Le Fèvre
ENS Lyon, France
Valentin.Le-Fevre@ens-lyon.fr

Louis-Claude Canon
FEMTO-ST, Besançon, France
Inria & ENS Lyon, France
Louis-Claude.Canon@univ-fcomte.fr

Yves Robert
ENS Lyon, France
Univ. Tennessee Knoxville, USA
Yves.Robert@inria.fr

Frédéric Vivien*
Inria & ENS Lyon, France
Frederic.Vivien@inria.fr

ABSTRACT

This work deals with scheduling and checkpointing strategies to execute scientific workflows on failure-prone large-scale platforms. To the best of our knowledge, this work is the first to target fail-stop errors for arbitrary workflows. Most previous work addresses soft errors, which corrupt the task being executed by a processor but do not cause the entire memory of that processor to be lost, contrarily to fail-stop errors. We revisit classical mapping heuristics such as HEFT and MINMIN and complement them with several checkpointing strategies. The objective is to derive an efficient trade-off between checkpointing every task (CKPTALL), which is an overkill when failures are rare events, and checkpointing no task (CKPTNONE), which induces dramatic re-execution overhead even when only a few failures strike during execution. Contrarily to previous work, our approach applies to arbitrary workflows, not just special classes of dependence graphs such as M-SPGs (Minimal Series-Parallel Graphs). Extensive experiments report significant gain over both CKPTALL and CKPTNONE, for a wide variety of workflows.

CCS CONCEPTS

• **Theory of computation** → **Scheduling algorithms**; • **Computer systems organization** → *Reliability*;

ACM Reference Format:

Li Han, Valentin Le Fèvre, Louis-Claude Canon, Yves Robert, and Frédéric Vivien. 2018. A Generic Approach to Scheduling and Checkpointing Workflows. In *ICPP 2018: 47th International Conference on Parallel Processing, August 13–16, 2018, Eugene, OR, USA*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3225058.3225145>

*All authors are with Laboratoire LIP, UMR 5668, Université de Lyon (CNRS, ENS Lyon, Inria, UCB Lyon), 69342 Lyon Cedex 07, France.

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

ICPP 2018, August 13–16, 2018, Eugene, OR, USA

© 2018 Copyright held by the owner/author(s). Publication rights licensed to the Association for Computing Machinery.

ACM ISBN 978-1-4503-6510-9/18/08...\$15.00

<https://doi.org/10.1145/3225058.3225145>

1 INTRODUCTION

This work deals with scheduling techniques to deploy scientific workflows on large parallel or distributed platforms. Scientific workflows are the archetype of HPC (High Performance Computing) applications, which are naturally partitioned into tasks that represent computational kernels. The tasks are partially ordered because the output of some tasks may be needed as input to some other tasks. Altogether, the application is structured as a DAG (Directed Acyclic Graph) whose nodes are the tasks and whose edges enforce the dependences. Nodes are weighted by the computational requirements (in flops) while edges are weighted by the size of communicated data (in bytes). Given a workflow and a platform, the problem of mapping the tasks onto the processors and to schedule them so as to minimize the total execution time, or makespan, has received considerable attention.

This classical mapping and scheduling problem has recently been revisited to account for the fact that errors and failures can strike during execution. Indeed, platform sizes have become so large that errors and failures are likely to strike at a high rate during application execution [8]. More precisely, the MTBF (Mean Time Between Failures) μ_P of the platform decreases linearly with the number of processors P , since $\mu_P = \frac{\mu_{\text{ind}}}{P}$, where μ_{ind} is the MTBF of each individual component (see Proposition 1.2 in [16]). Take $\mu_{\text{ind}} = 10$ years as an example. If $P = 10^5$ then $\mu_P \approx 50$ minutes and if $P = 10^6$ then $\mu_P \approx 5$ minutes: from the point of view of fault-tolerance, scale is the enemy.

Several approaches (see Section 6 for a review) have been proposed to mitigate the simplest instance of the problem, that of soft and silent errors. Soft errors cause a task execution to fail but without completely losing the data present in the processor memory. Local checkpointing (making a copy of all task input/output data) and task replication are the most widely used technique to address soft errors. Silent errors represent a different challenge than soft errors, in that they do not interrupt the execution of the task but corrupt its output data. However, their net effect is the same, since a task must be re-executed whenever a silent error is detected. A silent error detector is applied at the end of a task's execution, and the task must be re-executed from scratch in case of an error. Again, local checkpointing or replicating tasks and comparing outputs, are two common techniques to mitigate the impact of silent errors.

Fail-stop errors, or failures, are much more difficult to deal with. In the case of a fail-stop error (e.g., a crash due to a power loss) the execution of the processor stops, all the content of its memory is lost, and the computations have to be restarted from scratch, either on the same processor once it reboots or on a spare. The de-facto approach to handle such failures is Checkpoint/Restart (C/R), by which application state is saved to stable storage, such as a shared file system, throughout execution. The common strategy used in practice is *checkpoint everything*, or CKPTALL: all output data of each task is saved onto stable storage. For instance, in production Workflow Management Systems (WMSs) [1, 11, 24], the default behavior is that all output data is saved to files and all input data is read from files, which is exactly the CKPTALL strategy. While this strategy leads to fast restarts in case of failures, its downside is that it maximizes checkpointing overhead. At the other end of the spectrum would be a *checkpoint nothing* strategy, or CKPTNONE, by which all output data is kept in memory (up to memory capacity constraints) and no task is checkpointed. This corresponds to “in-situ” workflow executions, which has been proposed to reduce I/O overhead [26]. The downside is that, in case of a failure, a large number of tasks may have to be re-executed, leading to slow restarts. The objective of this work is to achieve a desirable trade-off between these two extremes. To the best of our knowledge, no general solution is available. We build upon our previous work [14] that was restricted to M-SPGs (Minimal Series-Parallel Graphs) [22]. In [14], we took advantage of the recursive structure of M-SPGs and used proportional mapping [18] for scheduling and checkpointing M-SPG workflows as sets of superchains. For general graphs, we have to resort to classical scheduling heuristics such as HEFT [20] and MINMIN [7], two reference scheduling algorithms widely used by the community. We provide extensions of HEFT and MINMIN that allow for a smaller subset of tasks to be checkpointed and lead to better makespans than the versions where each task (CKPTALL) or no task (CKPTNONE) is checkpointed.

The main contributions of this paper are the following:

- We deal with arbitrary dependence graphs, and require no graph transformation before applying our scheduling and checkpointing algorithms.
- We compare several mapping strategies and combine them with several checkpointing strategies.
- We design an event-based simulator to evaluate the makespan of the proposed solution. Indeed, computing the expected makespan of a solution is a difficult problem [14], and simple Monte-Carlo based simulations cannot be applied to general DAGs unless all tasks are checkpointed: otherwise, sampling the weight distribution for each task independently is not enough to compute the makespan, since a failure may involve re-executing several tasks.
- We report extensive experimental evaluation with both real-world and randomly generated workflows to quantify the performance gain achieved by the proposed approach.

The rest of the paper is organized as follows. First in Section 2, we work out an example to help understand the difficulty of the problem. Then we introduce the performance model in Section 3. We detail our scheduling and checkpointing algorithms in Section 4. We give experimental results in Section 5. Section 6 surveys the related work. Finally, we provide concluding remarks and directions for future work in Section 7.

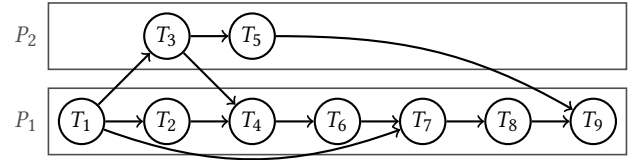


Figure 1: Schedule of a workflow with 9 tasks on 2 processors (each edge corresponds to a file dependence between tasks).

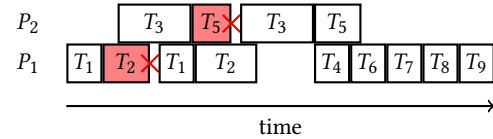


Figure 2: Sample execution of the workflow in Figure 1 without any checkpoint, with two failures striking during the execution of T_2 on P_1 and during that of T_5 on P_2 .

2 EXAMPLE

In this section, we illustrate the difficulty of deciding where to place checkpoints in a workflow. Consider the example of Figure 1 with 9 tasks, T_i , $1 \leq i \leq 9$, that have been mapped on 2 processors as shown on the figure. Note that this DAG cannot be reduced to an M-SPG and our previous approach [14] cannot be applied for this graph. While most tasks are assigned to processor P_1 , some tasks are assigned to the second processor, P_2 , to exploit the parallelism of the DAG. Any dependence between two tasks represents a file that is required to start the execution of the successor task; hence, $T_1 \rightarrow T_2$ represents a file produced by task T_1 that is required for the execution of task T_2 to start. Because T_1 and T_2 are both executed on processor P_1 , this file is kept in the memory of P_1 after T_1 completes. However, for the dependence $T_1 \rightarrow T_3$, because the tasks T_1 and T_3 are executed on different processors, the corresponding file must be retrieved by P_2 . Such a dependence between two tasks assigned to two different processors is called a *crossover dependence*.

In a first scenario, let us suppose that no task is checkpointed as showed in Figure 1: then if no failure strikes, the makespan will be the shortest possible, consisting only of the execution time of each task and of retrieving the necessary input files. However, as soon as a failure happens, we may need to restart the whole application from the very beginning. To study such a scenario, we need to explicit the memory management. Let us assume that once a processor has sent a file to another processor, then this file is deleted from the memory of the producing processor. For instance, as soon as P_2 has received from P_1 the file corresponding to the dependence $T_1 \rightarrow T_3$, this file is erased from the memory of P_1 . Remember that a failure wipes out the whole content of the memory of the struck processor. Thus, if a failure strikes during the execution of T_5 , to be able to re-attempt to execute T_5 , T_3 will need to be re-executed before (because the file $T_3 \rightarrow T_5$ is no longer available), which requires T_1 to be re-executed first (because the file $T_1 \rightarrow T_3$ is no longer available). Hence, a single failure in a part of the graph may require the re-execution of most of the workflow. Figure 2 shows an example of execution of the DAG when no task

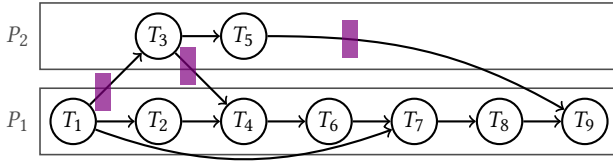


Figure 3: A purple *crossover* checkpoint is performed for each file produced by one processor and used by another one.

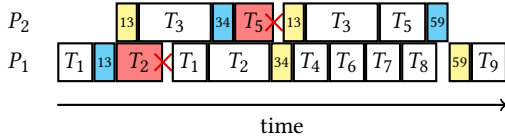


Figure 4: Sample execution of the application in Figure 3 with two failures striking during the execution of T_2 on P_1 and that of T_5 on P_2 , with crossover checkpoints. Label ij indicates the file from T_i to T_j . Now T_4 can start before the re-execution of T_3 since its output was checkpointed.

is checkpointed. To execute T_4 , we need both T_2 and T_3 to finish successfully, and that no fault strikes neither P_1 nor P_2 between the completion of these tasks and the start of T_4 . Here, T_2 does not finish so T_1 is re-executed. When P_2 fails, we need to re-execute T_3 , which requires input from T_1 . Luckily (!), P_1 already suffered from a failure, so T_1 has already been re-executed. Otherwise, we would have had to restart the execution of the whole workflow.

To avoid rolling back to the beginning in case of failures, we can try to place some checkpoints inside the workflow. As commonly assumed in workflow management systems [1, 11, 24], we do not rely on direct point-to-point communications between processors but instead assume that task input and output files are exchanged through the file system. Thus, any file produced by one processor and required by another processor is necessarily saved to, and then read from, stable storage. In the second scenario shown in Figure 3, we decide to checkpoint every crossover dependence (from T_1 to T_3 , T_3 to T_4 , and T_5 to T_9). An execution of that schedule is shown in Figure 4. Cyan boxes represent checkpoints while yellow boxes represent data being read. The transfer of file $T_1 \rightarrow T_3$ is done through a checkpointing phase on P_1 , followed by a reading phase on P_2 . We can see that thanks to the crossover checkpoints, T_4 does not need to wait for the completion of the second execution of T_3 anymore, as T_3 output data has already been checkpointed. Moreover, if only a failure on P_2 happened, instead of rolling back to task T_1 to re-execute T_3 as it was the case before, T_3 could have restarted directly (although the entire content of the processor memory is lost, so all inputs of T_3 must be recovered from stable storage after a downtime before the execution of T_3 can restart). The motivation to checkpoint all files involved in crossover dependences is to isolate the processors. Indeed, if all crossover files are checkpointed, a failure on a processor will never lead to the re-execution of a task successfully executed on another processor. Overall, we will

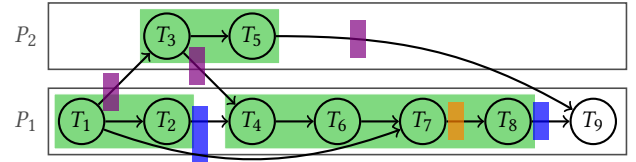


Figure 5: Blue *induced* checkpoints are used to isolate task sequences on a processor (labeled in green, such as the sequence T_4, T_6, T_7 and T_8 on P_1). Finally, additional checkpoints can be added inside an idle-free task sequence through a dynamic programming algorithm: the orange checkpoint corresponds to such an addition.

lose less time recomputing tasks or waiting for their second completion. However, reading from stable storage and checkpointing also take time. Finding the right trade-off is the main focus of this paper: deciding which tasks should be checkpointed, so that the overhead added by the checkpointing and reading of files is not more expensive than the re-execution of tasks.

We conclude by informally introducing examples of checkpointing strategies that achieve desirable trade-offs (see Section 4.2 for details). First, two additional checkpoints, in blue, called *induced* checkpoints, have been added in Figure 5. Their role is to secure the fast re-execution of tasks that are the target of a crossover dependence, namely T_4 and T_9 . The blue checkpoint after T_2 isolates the execution of the task sequence $S_1 = \{T_4, T_6, T_7, T_8\}$ on P_1 . To this purpose, it is necessary to checkpoint all intermediate results that may be used after the execution of T_2 : these are the files generated by previous tasks, namely $T_1 \rightarrow T_7$ and $T_2 \rightarrow T_4$. This way, when a failure strikes, previous tasks do not have to be restarted and the computation may be restarted directly from T_4 . This way, tasks in the sequence S_1 may be sequentially executed without idle time. It would not have been possible to include T_1 and T_2 in S_1 because T_4 could have waited for the completion of T_3 leading to idle time in some scenarios. Similarly, the second blue checkpoint isolates the execution of T_9 . Then, once the four tasks T_4, T_6, T_7 , and T_8 of the sequence S_1 have been “isolated” from other tasks, it is possible to use a dynamic programming algorithm similar to that used in [14] in order to introduce additional checkpoints. In Figure 5, a single additional checkpoint, in orange, is inserted after T_7 .

3 MODEL

This section details the execution and fault-tolerance models used to compare scheduling and checkpointing algorithms.

3.1 Execution Model

The execution model for a task workflow on a homogeneous system is represented as a Directed Acyclic Graph (DAG), $G = (V, E)$, where V is the set of nodes corresponding to the tasks, and E is the set of edges corresponding to the dependences between tasks. In a DAG, a node without any predecessor is called an *entry* node, while a node without any successor is an *exit* node. For a task T in G , $pred(T)$ and $succ(T)$ represent the set of its immediate predecessors and successors respectively. We say that a task T is *ready* if either it does not have any predecessor in the dependence graph, or if all its

predecessors have been executed. In this model, the execution time of a task $T_i \in V$ is w_i , i.e., its execution time in a failure-free execution. Each dependence $(T_i, T_j) \in E$ is associated with the cost $c_{i,j}$ to store/read the data onto/from stable storage. Before the execution of T_j on processor P_k , all input files needed by T_j must be present in the local memory of P_k and absent files must be *read* from the stable storage, which happens as late as possible. We ignore direct communication between processors because each data transfer between two processors (i.e., a *crossover dependence*) consists of writing to and reading from the stable storage. Alternatively, we also say that the file is checkpointed and then recovered.

3.2 Fault-Tolerance Model

In this work, each processor is a processing element that is subject to its own individual failures. Failures can strike a processor at any time, during either task execution or waiting time. Failure inter-arrival times are assumed to be Exponentially distributed. These failure-prone processors stop their execution once a failure strikes, i.e., we have fail-stop errors. When a fail-stop error strikes a processor, the whole content of its memory is lost and the computation it was performing must be restarted, either on the same processor after a reboot, or on a spare processor (e.g., taken from a pool of spare processors either specifically requested by the job submitter, or maintained by the resource management infrastructure).

Consider a single task T , with weight w , scheduled on such a processor, and whose input is stored on stable storage. It takes a time r to read that input data from stable storage, either for its first execution or after a failure. The total execution time W of T is a random variable, because several execution attempts may be needed before the task succeeds. We assume that failures are i.i.d. (independent and identically distributed) across the processors and that the failure inter-arrival times at each processor is Exponentially distributed with Mean Time Between Failures (MTBF) $\mu = 1/\lambda$. Let $\lambda \ll 1$ be the Exponential failure rate of the processor. With probability $e^{-\lambda(r+w)}$, no failure occurs, and W is equal to $r + w$. With probability $(1 - e^{-\lambda(r+w)})$, a failure occurs. For Exponentially distributed failures, the expected time to failure, knowing that a failure occurs during the task execution (i.e., in the next $r + w$ seconds), is $1/\lambda - (r+w)/(e^{\lambda(r+w)} - 1)$ [16]. After this failure, there is a downtime d , which is (an upper bound of) the time needed to reboot the processor or migrate to a spare. Then we start the execution again, first with the recovery r and then the work w . With a general model where an unbounded number of failures can occur during recovery and work, the expected time W to execute task T is given by $W = \left(\frac{1}{\lambda} + d\right) \left(e^{\lambda(r+w)} - 1\right)$ [16]. Now if the output data of task T is checkpointed, with a time c to write all of its output files onto stable storage, the total time becomes:

$$W = \left(\frac{1}{\lambda} + d\right) \left(e^{\lambda(r+w+c)} - 1\right). \quad (1)$$

Equation (1) assumes that failures can also occur during checkpoints, which is the most general model for failures. We also assume that failures may strike during the idle time (i.e., waiting time) of the processor (e.g., the power supply may fail). In the case of a sequence of non-checkpointed tasks to be executed on a processor P , the output data of each task resides in the memory of P for use

by subsequent tasks. When a failure strikes P , the entire memory content is lost and the whole task sequence must be re-executed.

4 SCHEDULING AND CHECKPOINTING ALGORITHMS

In this section, we first present heuristics to map tasks to processors. Then we propose three different checkpointing strategies that can be used simultaneously.

4.1 Scheduling heuristics

We map tasks to processors and schedule them using two classical scheduling heuristics, HEFT [20]¹ and MINMIN [7]. We run these heuristics as if the platforms were not subject to failures, that is, without considering checkpoints. Therefore, we decide first on which processor a task will be executed, and the order in which a processor will execute tasks, before deciding when and what to checkpoint (see Section 4.2). However, we present variants of HEFT and MINMIN, named HEFTC and MINMINC, that are specifically designed for our failure-prone framework.

Heterogeneous Earliest Finish Time first (HEFT) is presented as the HEFTC variant in Algorithm 1. The original HEFT algorithm comprises two phases. In a first *task prioritizing phase*, the bottom-level of all tasks is computed and tasks are ordered by non-increasing bottom-levels. The bottom-level of a task is the maximum length of any path starting at the task and ending in an exit task, considering that all communications take place [10]. In the second *processor selection phase*, the first unscheduled task is scheduled as early as possible on a processor that minimizes its completion time. In all cases, ties are broken arbitrarily. To these original two phases, we add a third one, the *chain mapping phase* (lines 7 and 8 of Algorithm 1). If the newly mapped task T is the head of a chain in the task graph, then this whole chain is mapped on the same processor as T , and the tasks will be executed consecutively. Ensuring that entire chain of tasks are scheduled on the same processor decreases the number of crossover dependences and thus, the time to checkpoint them. HEFTC has a complexity of $O(n^2)$ for a workflow with n tasks. During the processor selection phase, the earliest finish time of a task is computed in HEFTC while assuming that the newly mapped task must start *after* all tasks previously scheduled on that processor have completed. On the contrary, the original HEFT heuristic is allowed to perform backfilling following a classical insertion-based policy, as long as the completion time of no task is delayed. Allowing backfilling is more expensive at scheduling time but should lower the execution time (the complexity of HEFT with backfilling is also $O(n^2)$ with homogeneous processors). We do not allow backfilling for HEFTC because it could be antagonistic to the chain mapping phase if it led to backfill the head of the chain, but not the whole chain.

The MINMIN scheduling algorithm is presented in the MINMINC variant in Algorithm 2. The original MINMIN algorithm is a simple loop which, at each step, schedules the task that can finish the earliest among unscheduled tasks. Therefore, at each step it considers all ready tasks and, for each of them, all the processors. We (try to) improve this heuristic by adding a *chain mapping phase*

¹In fact, because we have homogeneous processors, we use MCP (Modified Critical Path) [25] with backfilling, which is exactly HEFT in this context.

Algorithm 1: HEFTC

```

1 Compute the bottom-level of all tasks by traversing the graph from
  the exit tasks
2 Sort the tasks by non-increasing values of their bottom-levels
3 while there are unscheduled tasks do
4   Select the first task  $T_i$ 
5    $k \leftarrow \arg \min_{1 \leq k \leq p} \text{EarliestFinishTime}(T_i, P_k)$ 
6   Schedule task  $T_i$  on processor  $P_k$ 
7   if  $T_i$  is the head of a chain of tasks then
8     Schedule the whole chain continuously on  $P_k$ 

```

exactly as previously (lines 5 and 6 of Algorithm 2). MINMINC has a complexity of $O(n^2p)$ for a workflow with n tasks and p processors.

Algorithm 2: MINMINC

```

1  $ReadyTasks \leftarrow \text{entry tasks}$ 
2 while there are unscheduled tasks do
3   Pick a task  $T \in ReadyTasks$  and a processor  $P$  such that the
  completion time of  $T$  on  $P$  is minimum among the Earliest
  Finish Times of all ready tasks
4   Schedule task  $T$  on processor  $P$ 
5   if  $T$  is the head of a chain of tasks then
6     Schedule the whole chain continuously on  $P$ 
7   Update  $ReadyTasks$ 

```

4.2 Checkpointing strategies

While the previous scheduling algorithms provide mappings of tasks to processors, it remains to decide which files must be checkpointed and when. This section introduces finer strategies than the two extremes solutions that consist of checkpointing no task or all tasks. These two extreme solutions, CKPTALL and CKPTNONE, are denoted with the suffixes NONE and ALL, respectively.

In principle, our model forbids direct communications between processors (see Section 3.1). However, for the sake of comparison, we make an exception for CKPTNONE: in the absence of any checkpoint with CKPTNONE, direct communications must be performed for each crossover dependence. We assume that, in this special case, transferring a file takes half the time needed to save it to and read it from stable storage. This special case is thus more efficient when files are large.

The minimum strategy that is required to avoid direct communications consists of checkpointing all files that must be transferred between any pair of processors, i.e., exactly the files corresponding to crossover dependences. Moreover, in this case, any failure on a processor will not require any re-execution on other processors. The strategy is denoted with a “C” in the checkpoint suffix.

For the next two additional strategies, we introduce a new type of checkpoints: *task checkpoints*. While a simple file checkpoint consists of writing to stable storage a file that corresponds to a dependence between two tasks, a task checkpoint consists of writing all files that (i) reside in memory on a processor; (ii) will be used later by tasks assigned to the same processor; and (iii) have not already been checkpointed. In the example in Section 2, for

each crossover dependence we did a simple file checkpoint rather than a full task checkpoint. A task checkpoint after task T_3 would have also checkpointed the file corresponding to the dependence $T_3 \rightarrow T_5$. A non-trivial task checkpoint for the example of Section 2 would be a task checkpoint for task T_2 . This checkpoint would require checkpointing the files corresponding to the dependences $T_2 \rightarrow T_4$ and $T_1 \rightarrow T_7$.

When a task checkpoint is performed after the execution of a task, multiple files may be checkpointed “at the same time” (either newly created files or previously created ones that will later be used). If several files are checkpointed, they are all checkpointed after the task completion, one after the other (in any order), and they can all be read again only when the last of them has been checkpointed. When absent from memory (following a failure or due to a crossover dependence), input files are read from stable storage as late as possible, just before the execution of the task that needs them. One could imagine optimizations where files (in a task checkpoint) would be checkpointed independently and as soon as possible, or in a carefully designed order. Such optimizations could lead to lower expected makespans in some cases. However, the interplay of file checkpoints and reads that could result from these optimizations may lead to slowdowns. This is the reason why we prefer our simpler scheme.

Checkpointing crossover dependences enable to isolate processors, in that there is no re-execution propagation from a processor to another. However, when a task is the target of a crossover dependence, its starting time is the maximum of the availability times of all its input files, and these files come from different processors. Therefore, its starting time may be delayed by failures occurring on other processors. Because failures can strike during idle time, it may be beneficial to try to use the potential waiting time by performing a task checkpoint of the task preceding the target task. This way, the whole content of the memory will be preserved, the cost of the checkpoint may be offset by some waiting time, and if a failure strikes during the remaining waiting time all input files remain available. Therefore, we propose a new checkpointing strategy denoted with “I” in the checkpoint suffix. This strategy consists of checkpointing all *induced* dependences. A dependence $T_i \rightarrow T_j$ is an *induced* dependence if T_i and T_j are scheduled on the same processor P and there exists a crossover dependence $T_k \rightarrow T_l$ such that T_l is scheduled on P after T_i and before T_j (or $T_l = T_j$). Checkpointing these induced dependences is done by performing a task checkpoint of the task preceding T_l on P . In the example of Section 2, the dependences $T_2 \rightarrow T_4$ and $T_1 \rightarrow T_7$ are both induced dependences because of the crossover dependence $T_3 \rightarrow T_4$.

So far, we have only introduced checkpoints to isolate processors, either to avoid failure propagation or to try to minimize the impact of processors having to wait from each other. We further consider checkpoints that more directly optimize expected total execution time. We present an additional strategy, denoted by the suffix “DP”, which adds additional checkpoints through a $O(n^2)$ dynamic programming algorithm, which is a transposition of that of [14]. This dynamic program considers a maximal sequence of consecutive tasks that are all assigned to the same processor, and that are isolated from other tasks: the sequence contains no checkpoint and none of its tasks is the target of a crossover dependence, except for its first task. Let T_1, \dots, T_k be such a sequence of tasks.

By definition, all input data produced by some previous tasks have been checkpointed. Then, the optimal expected time to execute this sequence is given by $Time(k)$ where $Time$ is defined as follows:

$$Time(j) = \min \left(T(1, j), \min_{1 \leq i < j} Time(i) + T(i + 1, j) \right)$$

where $T(i, j)$ is the expected time to execute tasks T_i to T_j provided that two task checkpoints are performed: one right before task T_i and one right after task T_j . Using the same reasoning as in Section 3.2, we can provide an upper bound on $T(i, j)$ as follows:

$$T(i, j) = \left(\frac{1}{\lambda} + d \right) \left(e^{\lambda(R_i^j + W_i^j + C_i^j)} - 1 \right)$$

where R_i^j (resp. W_i^j and C_i^j) is the sum of the recovery (resp. execution and checkpointing) costs of tasks T_i to T_j . The recovery costs concern all input files of these tasks that are on the stable storage, while the checkpointing costs concern all files that will be checkpointed when a task checkpoint is done after T_j . This is an upper bound, because when no failure strikes, some input files of tasks T_i to T_j may already be present in memory and will not be read from stable storage. Because we have no simple mean to know whether some failures had previously struck, we have to resort to this upper bound. This is a necessary condition to be able to reuse, in some way, the dynamic programming approach of in [14]. This algorithm requires, by construction, that induced dependences be checkpointed. However, we heuristically use it even when this condition is not satisfied. In this case, we take a maximal sequence while allowing tasks to be the target of crossover dependences, and behave as if these crossover dependences were not existing: we discard any potential waiting time that may be due to these crossover dependences (because we have no means to estimate them).

5 EXPERIMENTS

In this section, we describe the experiments conducted to assess the efficiency of the checkpointing strategies. In Subsection 5.1, we describe the parameters, applications and simulator used during our experimental campaign. We present our results in Subsection 5.2.

5.1 Experimental methodology

We consider workflows from real-world applications, namely representative workflow applications generated by the Pegasus Workflow Generator (PWG) [6], as well as the three most classical matrix decomposition algorithms (LU, QR, and Cholesky) [9], and randomly generated DAGs from the Standard Task Graph Set (STG) [19].

Pegasus workflows. PWG uses the information gathered from actual executions of scientific workflows as well as domain-specific knowledge of these workflows to generate representative and realistic synthetic workflows. We consider all five workflows [17] generated by PWG, including three M-SPGs (GENOME, LIGO, and MONTAGE) that are used to compare our new general approach with PROPCKPT, the strategy for M-SPGs proposed in [14]. We generate these workflows with 50, 300, and 700 tasks (these are the number of tasks given to the generator, the actual number of tasks in the generated workflows depend on the workflow shape). The task weights and file sizes are generated by PWG.

Matrix factorizations. We consider the three most classical factorizations of a $k \times k$ tiled matrix: LU, QR, and Cholesky factorizations. For each factorization, we perform experiments with $k = 6, 10,$ and $15,$ for a total of $3 \times 3 = 9$ DAGs with up to 1240 tasks. The number of vertices in the DAG depends on k as follows: the Cholesky DAG has $\frac{1}{3}k^3 + O(k^2)$ tasks, while the LU and QR DAGs have $\frac{2}{3}k^3 + O(k^2)$ tasks. There are 4 types of tasks in LU, QR, and Cholesky, which are labeled by the corresponding BLAS kernels [9], and their weights are based on actual kernel execution times as reported in [3].

Random graphs. The STG benchmark [19] includes 180 instances for each size of DAGs (from 50 to 5000). This set is often used in the literature to compare the performance of scheduling strategies. Instead of choosing part of the instances for each size, we did experiments on all instances of size 300 and 750.

Failure distribution. We consider different exponential processor failure rates. To allow for consistent comparisons of results across different DAGs (with different numbers of tasks and task weights), we fix the probability that a task fails, which we denote as p_{fail} , and then simulate the corresponding failure rate. Formally, for a given DAG $G = (V, E)$ and a given p_{fail} value, we compute the average task weight as $\bar{w} = \sum_{i \in V} w_i / |V|$, where w_i is the weight of the i -th task in V . We then pick the failure rate λ such that $p_{\text{fail}} = 1 - e^{-\lambda \bar{w}}$. We conduct experiments for three p_{fail} values: 0.01, 0.001, and 0.0001.

Checkpointing costs. An important factor that influences the performance of checkpointing strategies, and more precisely of the checkpointing and recovery overheads, is the data-intensiveness of the application. We define the Communication-to-Computation Ratio (CCR) as the time needed to store all the files handled by a workflow (input, output, and intermediate files) divided by the time needed to perform all the computations of that workflow on a single processor. For Pegasus workflows, LU, QR, and Cholesky, we vary the CCR by scaling file sizes by a factor. As STG only provides task weights, we compute the average communication cost as $\bar{c} = \bar{w} \times CCR$. Communication costs are generated with a lognormal distribution with parameters $\mu = \log(\bar{c}) - 2$ and $\sigma = 2$ to ensure an expected value of \bar{c} . This distribution with parameter $\sigma = 2$ has been advocated to model file sizes [12]. This allows considering and quantifying the data-intensiveness of all workflows in a coherent manner across experiments and workflow classes and configurations.

Reference strategies. We compare our strategies to the two extreme approaches CKPTALL and CKPTNONE. For each parameter setting of each workflow, we run 10,000 random simulations and approximate the makespan by the observed average makespan.

Simulator. We implemented a discrete event simulator. The C++ code is available at <http://github.com/vlefevre/task-graph-simulation>. A full description of the simulator can be found in [15]. This simulator computes the number of file checkpoints taken, the number of task checkpoints taken, the number of failures, the total time spent checkpointing data, and the execution time.

5.2 Results

Due to space limitations, we only report here a subset of our simulations results. For instance, we only report on the Cholesky factorization as results are similar for LU and QR, and on CYBERSHAKE

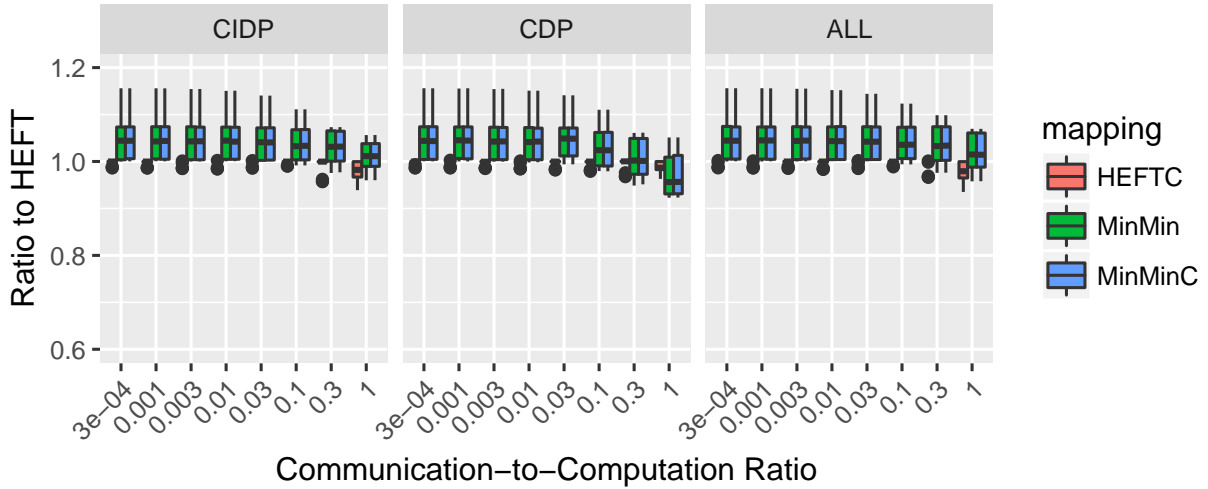


Figure 6: Relative performance of the four task mapping and scheduling strategies for Cholesky.

and SIPHT as results are similar for GENOME, LIGO, and MONTAGE. All results can be found in the companion research report [15].

First we compare the four considered task mapping and scheduling strategies: HEFT and MINMIN, with their chain-mapping variants HEFTC and MINMINC. Figure 6 presents such a comparison for Cholesky using boxplots². On this figure, the lower the better and the baseline at 1 is the performance of HEFT. The chain-mapping variants have the same performance or improve that of their basic counterparts, especially when communications are expensive (rightmost parts of the graphs). The other conclusion is that MINMIN (resp. MINMINC) almost always achieves same or worse performance than HEFT (resp. HEFTC). This is explained by the fact that HEFT and HEFTC take into account the critical path of workflows. These trends are representative of the trends observed for all considered graphs and workflows (see [15] for the other graphs) but suffer from some exceptions. The chain-mapping variants can be superceded by their basic counterparts for workflows that do not include any chains (like LU), because the basic variants use backfilling. However, backfilling sometimes backfires, even in the absence of chains, like for SiphT where HEFTC can decrease the expected makespan by more than 30% with respect to HEFT (see [15]). Overall, of the four considered task mapping and scheduling heuristics, HEFTC never achieves significantly bad performance, and most of the time achieves the best performance. This is the reason why we focus on it in the remainder of this section.

Figures 7 through 9 present the expected makespans achieved by CDP, CIDP³, and NONE divided by that of ALL when the Communication-to-Computation Ratio increases. The lower the better and data points below the $y = 1$ line denote cases in which these strategies outperform the competitor ALL. Each figure shows results for workflows with different number of tasks (each line of subfigure is for a different size, the number of tasks being reported on the rightmost column), for various number of processors P (different

line styles), and for the three p_{fail} values (0.0001, 0.001, 0.01). On these figures, we report in black above the horizontal axis the average number of failures that occur for the 10,000 random trials for each setting. The other two lines of numbers are the number of checkpointed tasks for CDP and CIDP, each number is printed with the same color as the curve of the corresponding strategy.

CIDP never achieves worse performance than ALL: either it achieves a similar performance or it outperforms ALL (for Montage, Genome and Ligo, see [15]), especially when communications, and thus checkpoints, are expensive. When checkpoints come for free, ALL and CIDP have the same performance as they both checkpoint all tasks.

In the majority of cases, CDP also achieves similar or better performance than ALL. As explained in Section 4, the dynamic programming algorithm is well-defined for CIDP, which checkpoints all induced dependences. However, CDP tries to save some checkpointing overhead by not systematically checkpointing induced dependences. As a consequence, the dynamic programming algorithm estimations of expected execution times may be inaccurate, which explains the sometimes bad performance of CDP. There are only a couple of CCR values for CYBERSHAKE for which CDP achieves a significant worse performance than ALL. On the contrary, CDP often has better performance than CIDP when checkpointing cost is high. CDP never checkpoints more tasks than CIDP. Depending on the checkpointing cost and failure rate, CDP can lead to significant improvement over ALL. For workflows as dense as LU, we save more than 10% when $CCR = 1$ for both strategies (see[15]), and CDP even achieves 35% saving for SIPHT. As the CCR decreases, the ratio converges to 1 because all strategies checkpoint all tasks.

CDP and CIDP achieve better results than NONE except when (i) checkpoints are expensive (high CCR) and/or (ii) failures are rare (low p_{fail}). In these cases, checkpointing is a losing proposition, and yet our strategies, by design, always checkpoints some files (all crossover files and even induced dependences for CIDP). The optimal approach is then to bet that no failure will happen and to

²Each boxplot consists of a bold line for the median, a box for the quartiles, whiskers that extend at most to 1.5 times the interquartile range from the box and additional points for outliers.

³CDP and CIDP correspond to the designations introduced in Section 4.2.

restart the whole workflow execution from scratch upon the very rare occurrence of a failure. NONE becomes worse whenever there are more failing tasks, i.e., when the failure rate increases, and/or when the number of tasks increases. When the failure rate is high and the workflows are large, the relative expected makespan of NONE is so high that it does not appear in the plots.

Figure 10 presents the aggregated results for the 180 STG random DAGs with boxplots. The trends on these graphs are the same as already reported. This confirms the generality of our conclusions.

Finally, we compare our new general approach with PROPCKPT, the approach specific to M-SPGs that we proposed in [14]. Figure 11 presents this comparison for Montage, which is the worst case for our approaches compared to Genome and Ligo. Overall, the new approaches perform better than PROPCKPT.

6 RELATED WORK

Checkpointing workflows has received considerable attention in the recent years, but no satisfactory solution has yet been proposed for fail-stop failures and general DAGs.

Many authors have considered soft errors, by which a task execution fails but does not lead to completely losing the data present in the processor memory. Fail-stop errors have far more drastic consequences than soft errors as they induce the loss of all data present in memory. Therefore they require different solutions. As discussed in Section 1, silent errors do not interrupt the execution of the task but corrupt its output data. Their net effect is the same, since a task must be re-executed whenever a silent error is detected. Their detection requires the use of some silent error detectors at the end of a task's execution. As we only consider fail-stop errors we do not need to use fault detectors. See our previous work [14] for an overview of the related work on soft and silent errors.

Relatively few published works have studied fail-stop failures in the context of workflow applications. When the workflow is a linear chain of tasks, the problem of finding the optimal checkpoint strategy has been solved by Toueg and Babaoglu [21] using a dynamic programming algorithm. This algorithm was later extended in [5] to cope with both fail-stop and silent errors simultaneously. When the workflow is general but comprised of parallel tasks that each executes on the whole platform, the problem of placing checkpoints is NP-complete for simple join graphs [4]. Existing work in the most general context, i.e., when tasks of a workflow do not necessarily span the whole platform, diverges from ours as follows: either there is a limit to the number of failures that an execution can cope with [23], or the optimization objective is reliability [2], meaning that the application execution can fail altogether. The only exception that we are aware of is our previous work [14]. The limitation of that work was different: the proposed solution could only deal with workflows whose structure was a Minimal Series-Parallel Graph (a generalization of Series-Parallel Graph).

To the best of our knowledge, this work is the first approach (beyond application-specific solutions) that (i) does not resort to linearizing the entire workflow as a chain of (parallel) tasks; (ii) can be applied to any workflow; (iii) can cope with an arbitrary number of failures; (iv) always guarantees a successful application execution; and (v) minimizes the (expectation of) the application

execution time. As a result, we propose the first DAG scheduling/checkpointing algorithm that allows arbitrary workflows to execute concurrently on multiple failure-prone processors in standard task-parallel fashion.

7 CONCLUSION

This work tackles the challenging problem of executing arbitrary workflows on homogeneous processors, with reasonable performance in presence of failures but without incurring a prohibitive cost when no failure strikes. While CKPTALL meets the first objective by expensively checkpointing every task and CKPTNONE meets the second one by avoiding any checkpoint at all, we propose new strategies that provide different trade-offs between these two extremes. First, all crossover dependences, corresponding to file transfers between processors, are checkpointed, which prevents re-execution propagation between processors in case of failure. Then, a DP (Dynamic Programming) solution is used to insert additional checkpoints to minimize the expected completion time. Additional (induced) checkpoints may be added prior to the DP execution to provide it with more accurate information. Moreover, different mapping strategies that extend classical ones to reduce the number of checkpoints were also proposed. To the best of our knowledge, these new strategies are the first to be tuned to minimize the need for checkpointing while mapping tasks. Extensive experiments with a discrete event simulator, conducted for both synthetic and realistic instances, show that our approaches significantly outperform CKPTALL and CKPTNONE in most scenarios.

Future work will aim at extending our approach to workflows with parallel moldable tasks [13]. Such an extension raises yet another significant challenge: now the number of processors assigned to each task becomes a parameter to the proposed solutions, with a dramatic impact on both performance and resilience.

ACKNOWLEDGMENTS

This work is partially supported by the LABEX MILYON (ANR-10-LABX-0070), being operated by the French National Research Agency (ANR), and by a JORISS grant.

REFERENCES

- [1] Michael Albrecht, Patrick Donnelly, Peter Bui, and Douglas Thain. 2012. Makeflow: A portable abstraction for data intensive computing on clusters, clouds, and grids. In *1st ACM SIGMOD Workshop on Scalable Workflow Execution Engines and Technologies*. ACM.
- [2] Ismail Assayad, Alain Girault, and Hamoudi Kalla. 2004. A Bi-Criteria Scheduling Heuristic for Distributed Embedded Systems under Reliability and Real-Time Constraints. In *Dependable Systems Networks (DSN)*. IEEE.
- [3] Cédric Augonnet, Samuel Thibault, Raymond Namyst, and Pierre-André Wacrenier. 2011. StarPU: a unified platform for task scheduling on heterogeneous multicore architectures. *Concurrency and Computation: Practice and Experience* 23, 2 (2011), 187–198.
- [4] Guillaume Aupy, Anne Benoit, Henri Casanova, and Yves Robert. 2016. Scheduling computational workflows on failure-prone platforms. *Int. J. of Networking and Computing* 6, 1 (2016), 2–26.
- [5] Anne Benoit, Aurélien Cavelan, Yves Robert, and Hongyang Sun. 2016. Assessing general-purpose algorithms to cope with fail-stop and silent errors. *ACM Trans. Parallel Computing* 3, 2 (2016).
- [6] Shishir Bharathi, Ann Chervenak, Ewa Deelman, Gaurang Mehta, Mei-Hui Su, and Karan Vahi. 2008. Characterization of scientific workflows. In *Workflows in Support of Large-Scale Science (WORKS)*. IEEE, 1–10.
- [7] Tracy D Braun, Howard Jay Siegel, Noah Beck, Ladislau L Bölöni, Muthucumaru Maheswaran, Albert I Reuther, James P Robertson, Mitchell D Theys, Bin Yao, Debra Hensgen, et al. 2001. A comparison of eleven static heuristics for mapping

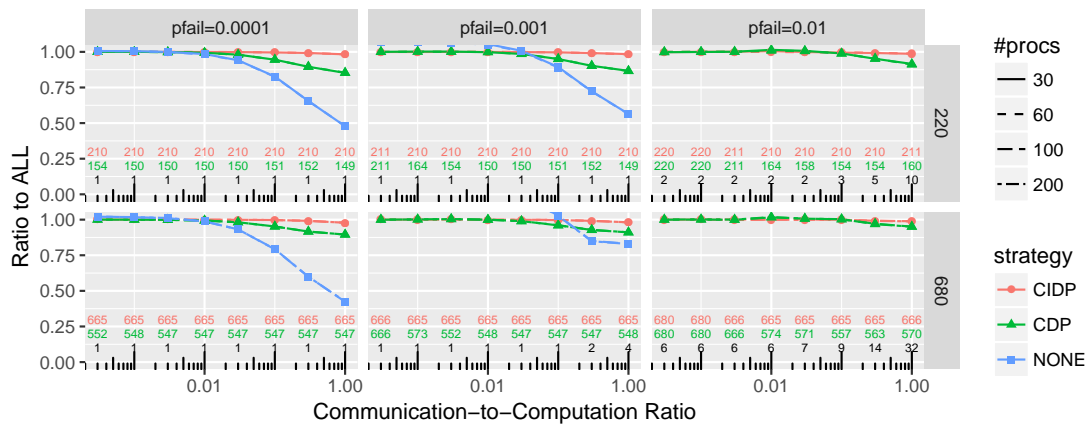


Figure 7: Performance of the checkpointing strategies for Cholesky using HEFTC for task mapping and scheduling.

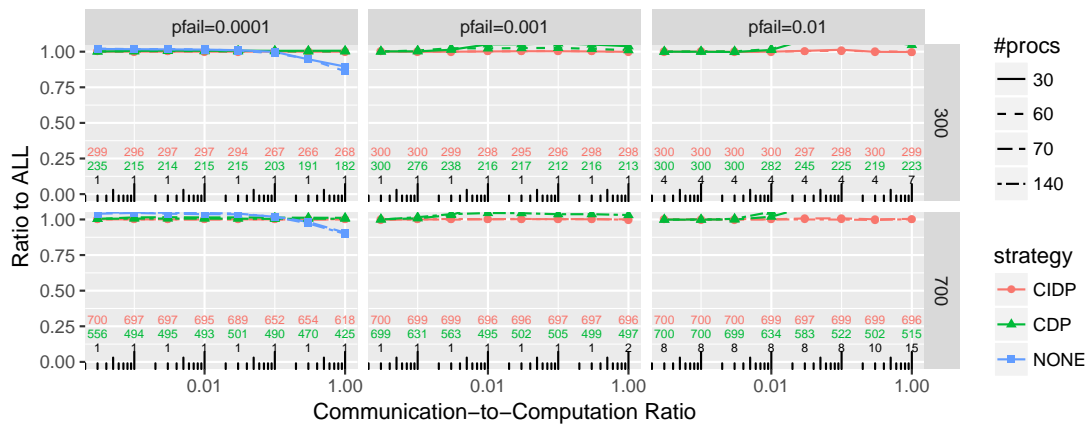


Figure 8: Performance of the checkpointing strategies for CyberShake using HEFTC for task mapping and scheduling.

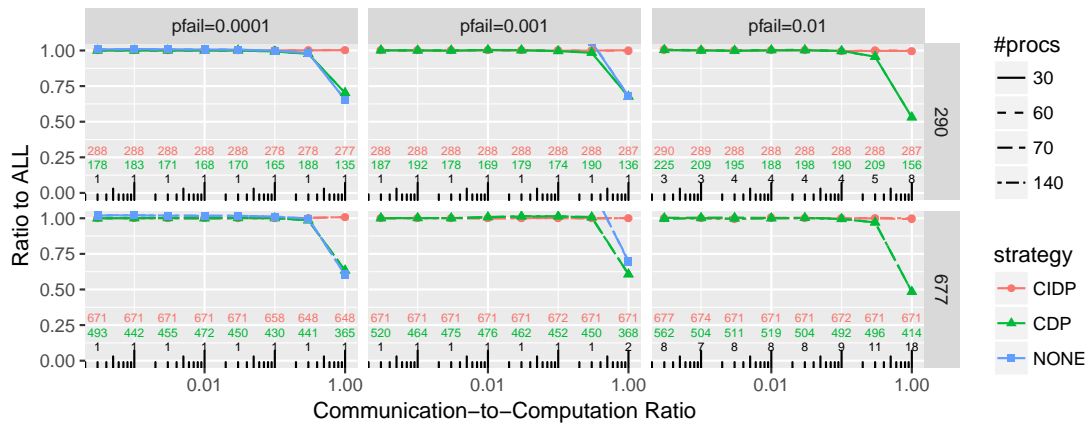


Figure 9: Performance of the checkpointing strategies for Siptt using HEFTC for task mapping and scheduling.

a class of independent tasks onto heterogeneous distributed computing systems. *Journal of Parallel and Distributed computing* 61, 6 (2001), 810–837.

[8] Franck Cappello, Al Geist, William Gropp, Sanjay Kale, Bill Kramer, and Marc Snir. 2014. Toward Exascale Resilience: 2014 update. *Supercomputing frontiers and innovations* 1, 1 (2014).

[9] Jaeyoung Choi, Jack J Dongarra, L Susan Ostrouchov, Antoine P Petit, David W Walker, and R Clint Whaley. 1996. Design and implementation of the ScaLAPACK LU, QR, and Cholesky factorization routines. *Scientific Programming* 5, 3 (1996), 173–184.

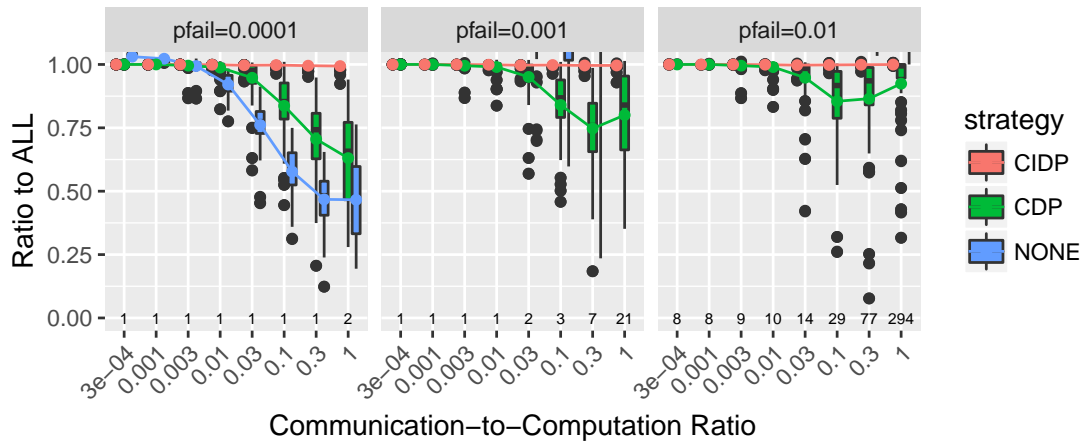


Figure 10: Average performance of the checkpointing strategies for the STG task graphs using HEFTC for task mapping and scheduling.

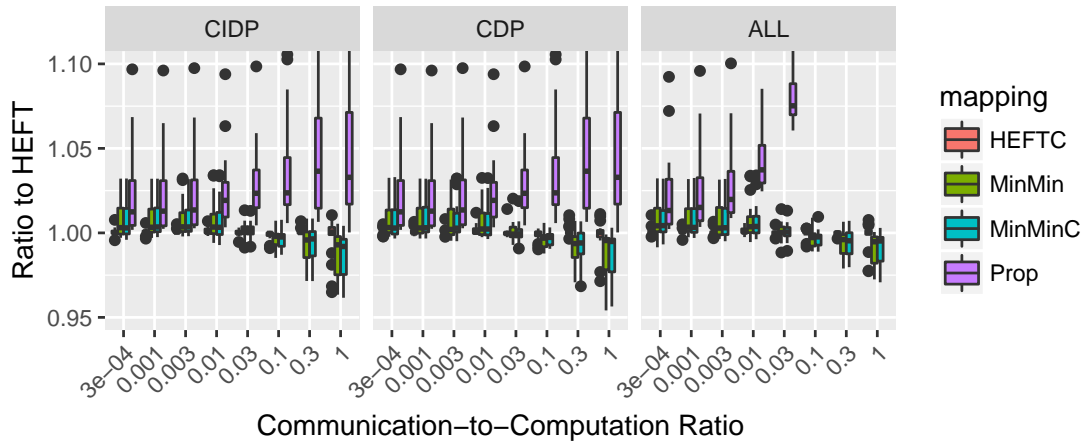


Figure 11: Relative performance of the four task mapping and scheduling strategies and of PROPCKPT for Montage.

- [10] Alain Darte, Yves Robert, and Frédéric Vivien. 2000. *Scheduling and automatic parallelization*. Birkhäuser.
- [11] Ewa Deelman, Karan Vahi, Gideon Juve, Mats Rynge, Scott Callaghan, Phil J. Maechling, Rajiv Mayani, Weiwei Chen, Rafael Ferreira da Silva, Miron Livny, and Kent Wenger. 2015. Pegasus, a Workflow Management System for Science Automation. *Future Generation Computer Systems* 46 (2015), 17–35.
- [12] Allen B Downey. 2001. The structural cause of file size distributions. In *Modeling, Analysis and Simulation of Computer and Telecommunication Systems, 2001. Proceedings. Ninth International Symposium on*. IEEE, 361–370.
- [13] Maciej Drozdowski. 2009. *Scheduling for Parallel Processing*. Springer.
- [14] Li Han, Louis-Claude Canon, Henri Casanova, Yves Robert, and Frédéric Vivien. 2018. Checkpointing workflows for fail-stop errors. *IEEE Trans. Comput.* (2018).
- [15] Li Han, Valentin Le Fèvre, Louis-Claude Canon, Yves Robert, and Frédéric Vivien. 2018. A generic approach to scheduling and checkpointing workflows. Research report 9167. INRIA.
- [16] Thomas Héroult and Yves Robert (Eds.). 2015. *Fault-Tolerance Techniques for High-Performance Computing*. Springer Verlag.
- [17] Pegasus. 2014. Pegasus Workflow Generator. <https://confluence.pegasus.isi.edu/display/pegasus/WorkflowGenerator>. (2014).
- [18] Alex Pothen and Chunguang Sun. 1993. A mapping algorithm for parallel sparse Cholesky factorization. *SIAM J. on Scientific Computing* 14, 5 (1993), 1253–1257.
- [19] Takao Tobita and Hironori Kasahara. 2002. A standard task graph set for fair evaluation of multiprocessor scheduling algorithms. *Journal of Scheduling* 5, 5 (2002), 379–394.
- [20] Haluk Topcuoglu, Salim Hariri, and Min-you Wu. 2002. Performance-effective and low-complexity task scheduling for heterogeneous computing. *IEEE transactions on parallel and distributed systems* 13, 3 (2002), 260–274.
- [21] Sam Toueg and Özalp Babaoglu. 1984. On the Optimum Checkpoint Selection Problem. *SIAM J. Comput.* 13, 3 (1984).
- [22] Jacobo Valdes, Robert E. Tarjan, and Eugene L. Lawler. 1979. The Recognition of Series Parallel Digraphs. In *Proc. of STOC'79*. ACM, 1–12.
- [23] P. Wang, K. Zhang, R. Chen, H. Chen, and H. Guan. 2014. Replication-Based Fault-Tolerance for Large-Scale Graph Processing. In *2014 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*. 562–573.
- [24] Katherine Wolstencroft, Robert Haines, Donal Fellows, Alan Williams, David Withers, Stuart Owen, Stian Soiland-Reyes, Ian Dunlop, Aleksandra Nenadic, Paul Fisher, et al. 2013. The Taverna workflow suite: designing and executing workflows of Web Services on the desktop, web or in the cloud. *Nucleic acids research* (2013), gkt328.
- [25] M. Y. Wu and D. D. Gajski. 1990. Hypertool: a programming aid for message-passing systems. *IEEE Trans. Parallel Distributed Systems* 1, 3 (1990), 330–343.
- [26] Fan Zhang, Ciprian Docan, Manish Parashar, Scott Klasky, Norbert Podhorszki, and Hasan Abbasi. 2012. Enabling In-situ Execution of Coupled Scientific Workflow on Multi-core Platform. In *Proc. 26th IEEE IPDPS*. 1352–1363.