

# Computing the expected makespan of task graphs in the presence of silent errors

Henri Casanova\*, Julien Herrmann† and Yves Robert†‡

\*Information and Computer Sciences Dept., University of Hawai‘i, USA

Email: henric@hawaii.edu

†École Normale Supérieure de Lyon, France

Email: {julien.herrmann|yves.robert}@ens-lyon.fr

‡University of Knoxville, USA

**Abstract**—Applications structured as Directed Acyclic Graphs (DAGs) of tasks correspond to a general model of parallel computation and thus arise in most domains (often as “scientific workflows”). A key question that is the subject of large literature is that of DAG scheduling, for which several list-scheduling heuristics have been proposed and shown to be effective in practice. Many of these heuristics make scheduling decisions based on path lengths in the DAG. At large scale, however, compute platforms and thus tasks are subject to various types of failures with no longer negligible probabilities of occurrence. Failures that have recently received increasing attention are “silent errors,” which cause a task to produce incorrect results even though it ran to completion. Tolerating silent errors is done by checking the validity of the results and re-executing the task from scratch in case of an invalid result. The execution time of a task then becomes a random variable, and so are path lengths. Unfortunately, computing the expected makespan of a DAG (and equivalently computing expected path lengths in a DAG) is a computationally difficult problem. Consequently, designing effective scheduling heuristics is preconditioned on computing accurate approximations of the expected makespan. In this work we propose an algorithm that computes a first order approximation of the expected makespan of a DAG when tasks are subject to silent errors. We compare our proposed approximation to previously proposed such approximations for three classes of application graphs from the field of numerical linear algebra. Our evaluations quantify approximation error with respect to a ground truth computed via a brute-force Monte Carlo method. We find that our proposed approximation outperforms previously proposed approaches, and in particular leads to large reductions in approximation error for low (and realistic) failure rates.

## I. INTRODUCTION

This paper introduces a new algorithm to approximate the expected makespan of a workflow application, i.e.,

an application structured as a Directed Acyclic Graph (DAG) of tasks, in which tasks can fail and must be re-executed. A key question when executing workflows on parallel platforms is the scheduling of tasks on the available compute resources, or processors. When not considering task failures, list scheduling algorithms are the de-facto standard [1], and tools are available that use such algorithms for scheduling workflow applications onto large-scale platforms in practice [2], [3]. The prevalent list scheduling heuristics prioritize tasks with large bottom levels. The bottom-level of a task is defined as the longest path from that task to the end of the execution, assuming unlimited processors. This heuristic is known as CP-scheduling (*Critical Path* scheduling [4], [5], [6]), and it has been extended to handle heterogeneous environments (see the HEFT algorithm [7]).

It is well-documented that large-scale platforms are increasingly subject to errors that can cause task failures. In particular, the occurrence of “silent errors” (or SDCs, for Silent Data Corruptions) has been recently identified as one of the major challenges for Exascale [8]. Silent errors can be caused by external causes, including cosmic radiations or packaging pollution. In addition, silent errors can occur when using DVFS (Dynamic Voltage Frequency Scaling) to reduce energy consumption. For instance, at low voltage the probability that a task produces incorrect results increases [9], [10].

Regardless of their causes, an effective approach for avoiding propagating results corrupted by silent errors is to use a verification mechanism after executing each task [11]. When an error is detected, the task is then re-executed. The verification mechanism can be general-purpose (e.g., based on replication [12], with re-execution only when the two outputs do not match) or application-specific. Many application-specific error detection methods are available for classical High perfor-

mance Computing (HPC) applications, such as approximate re-execution for ODE and PDE solvers [13], orthogonality checks for Krylov-based sparse solvers [14], [15]), or Algorithm-Based Fault Tolerance (ABFT) for numerical linear algebra [16].

Silent errors make it challenging to define efficient list scheduling algorithms to schedule workflow applications. After the first execution of a task, the error detector is used to check the result. If the result is correct, the task’s execution is marked as successful, and its successor tasks are marked as ready. If the result is incorrect, then the task must be executed (and, given that the first execution has failed, the second execution will typically succeed with very high probability). While this scheme is conceptually simple, it greatly complicates the computation of the bottom-level of a task. And yet, computing the expected length of the longest path in a DAG with unlimited processors (or equivalently the expected bottom-level of a task in the DAG), is key to designing silent-error-aware versions of effective list scheduling heuristics (CP-scheduling, HEFT).

It is known that computing the expected length of the longest path in a DAG whose task weights are probabilistic is a difficult problem [6]. Even in the case in which each task is re-executed at most once, i.e., when task weights are random variables taking only two discrete values, the problem remains #P-complete [17] (see Section II for a detailed discussion). In this work we develop an algorithm to compute an accurate first-order approximation of the expected length of the longest path in a general DAG in which tasks are subject to silent errors, and whose execution lengths can take two different values, depending upon there is a re-execution or not. More specifically, our contributions are:

- We develop an exact first-order approximation of the expected makespan of a general DAG, which can be computed in polynomial time;
- We compare our approach to two previously proposed approximations for three classes of DAGs from numerical linear algebra computations;
- We quantify approximation errors via comparison to a brute-force Monte Carlo approach; and
- We show that our proposed approximation leads to lower or similar error than previously proposed approximations, and importantly to much lower error when the probability of task failure is low.

This paper is organized as follows. Section II reviews related work. Section III formalizes the problem that we address and states assumptions. Section IV describes our

proposed approximation. Section V presents evaluation results. Section VI concludes with a brief summary of results and with perspectives on future work.

## II. RELATED WORK

In this section we first review previous works that focus on scheduling algorithms for probabilistic DAGs. We then review the relevant literature on silent errors.

### A. Expected makespan of probabilistic 2-state DAGs

Computing the expected makespan of a DAG whose task execution times obey arbitrary probability distributions is known to be a difficult problem, even with unlimited processors. When task weights are fixed (deterministic), the makespan is the length of the longest path in the graph (also called critical path). But assume instead a *probabilistic 2-state DAG* where task weights obey a simple 2-state probability distribution: task  $T_i$  has weight  $a_{i,1}$  with probability  $p_i$  and weight  $a_{i,2}$  with probability  $1 - p_i$ . Assume also that all the probability distributions of all tasks are independent. The makespan of the DAG is a random variable. It is known that computing its probability distribution, or even just its expected value, is a #P-complete problem [17]. Recall that the class of #P problems is the class of counting problems corresponding to NP decision problems [18], [19], [20], and that #P-complete problems are at least as hard as NP-complete problems.

An informal explanation of why computing the expected makespan of probabilistic 2-state DAGs is a difficult combinatorial problem, is as follows. The main intuitive reason is that the expected value of the maximum of two random variables is not the maximum of their expectations. As a result, when computing the length of a path in the DAG, one must keep track of all possible values for the starting time of each task, together with their probabilities, and there may be an exponential number of such values [21].

In practice, there are three standard methods to compute the expected makespan of a probabilistic 2-state DAG, as described hereafter.

1) *Monte Carlo simulations*: The Monte Carlo approach works as follows [22], [23]: For each task in the DAG, a value of its weight is sampled from its probability distribution. Once this is done, the DAG is deterministic and its longest path can be computed as explained in Section III. One then repeats this operation for a large number of trials, generating a new value at each trial. The set of these values empirically approaches

the actual distribution of the DAG makespan as the number of trials increases.

An interesting question is that of determining the number of trials to obtain a high confidence level in the result. We refer the reader to the relevant discussion in [24]. A key drawback of the Monte Carlo approach is that it is compute-intensive since the necessary number of trials is typically high. In this work, we only use Monte Carlo as a ground truth to assess the accuracy of our and previously proposed algorithms that compute approximations of the expected makespan. Hence, instead of determining a minimal number of trials, we use a very large number of trials so as to guarantee that our ground truth is accurate.

2) *Approximation by a series-parallel graph*: Basic probability theory tells us how to compute the probability distribution of the sum of two random variables (by a convolution) and of the maximum of two random variables (by taking the product of their cumulative density functions). This simple consideration leads to an exact method to compute the expected makespan when the DAG is series-parallel (see [25] for a definition of series-parallel graphs). However, the problem with probabilistic 2-state series-parallel graph remains NP-complete in the weak sense and admits a pseudo-polynomial solution [21].

When the DAG is not series-parallel, one approach is to approximate it by a series-parallel graph, which is constructed iteratively, first by a sequence of reductions and then by duplicating some vertices. Dodin’s method [26] constructs such an approximated series-parallel graph, whose expected makespan is used to estimate that of the original DAG. See [21], [24] for a detailed description of Dodin’s method. We include this method in our quantitative experiments in Section V.

3) *Approximation with normality assumption*: The central-limit theorem states that the sum of independent random variables tends to be normally distributed as the number of variables increases. The expected makespan of the DAG is a combination of sums and maximums of the original task weights, so a popular approach proposed by Sculli [27] is based on the *normality assumption*:

- Approximate the distribution of each task by a normal distribution of same mean and variance. This step has constant cost per task for probabilistic 2-state DAGs.
- Use Clarke’s formulas in [28] to compute the mean and variance of the sum and maximum of two

(correlated) normal distributions, and then assuming that they also follow normal distributions.

- Traverse the original DAG and compute the mean and variance of the makespan.

See [24] for a full description of Sculli’s method, which we evaluate experimentally in Section V.

### B. Silent errors

Considerable efforts have been directed at verification techniques to handle silent errors. A guaranteed, general-purpose verification is only achievable with expensive techniques, such as process replication [12], [29] or redundancy [30], [31]. However, application-specific information can be exploited to decrease the verification cost. Algorithm-based fault tolerance (ABFT) [32], [16], [33] is a well-known technique to detect errors in linear algebra kernels using checksums. Various techniques have been proposed in other application domains. Benson et al. [13] compare a higher-order scheme with a lower-order one to detect errors in the numerical analysis of ODEs. Sao and Vuduc [15] investigate self-stabilizing corrections after error detection in the conjugate gradient method. Heroux and Hoemmen [34] design a fault-tolerant GMRES capable of converging despite silent errors. Bronevetsky and de Supinski [35] provide a comparative study of detection costs for iterative methods.

Recently, detectors based on data analytics have been proposed to serve as partial verifications [36], [37], [38]. These detectors use interpolation techniques, such as time series prediction and spatial multivariate interpolation, on scientific dataset to offer large error coverage at the expense of a negligible overhead. Although not perfect, the accuracy-to-cost ratios of these techniques tend to be very high, which makes them attractive alternatives at large scale.

As mentioned in Section I, lowering the voltage/frequency is believed to have an adverse effect on system reliability [9], [10]. In particular, many papers (e.g., [10], [39], [40], [41]) have assumed the following exponential error rate model:

$$\lambda(s) = \lambda_0 \cdot 10^{\frac{d(s_{max}-s)}{s_{max}-s_{min}}}, \quad (1)$$

where  $\lambda_0$  denotes the average error rate at the maximum speed  $s_{max}$ ,  $d > 0$  is a constant indicating the sensitivity of error rate to voltage/frequency scaling, and  $s_{min}$  is the minimum speed. This model suggests that the error rate increases exponentially with decreased processing speed, which is a result of decreasing the voltage/frequency and hence lowering the circuit’s critical charge (i.e., the minimum charge required to cause an error in the circuit).

This is bad news for resilience, because these studies suggest that minimizing energy via DVFS techniques will also lead to an increased number of silent errors.

After detecting a silent error via some verification mechanism, the task is re-executed a second time. Because silent errors are not detected when they occur, but only at the end of the execution of the task (by the verification mechanism), the task must be fully re-executed, and its weight thus doubles. This is the reason why we consider probabilistic 2-state DAGs where the weight of a task  $T$  is its initial cost  $a$  if there is no error, and  $2a$  otherwise.

### III. PROBLEM STATEMENT

We consider a general model of computation in which an application is structured as a Directed Acyclic Graph, in which vertices represent tasks and edges represent task precedence. More formally, let  $G = (V, E)$  be a DAG, with  $V$  a set of tasks, and  $E \subset V \times V$  a set of edges. For each task  $i$ , let  $a_i$  be its weight, i.e., its failure-free execution time. For each task  $i \in V$ , let  $Pred(i) = \{j \in V \mid (j, i) \in E\}$  and  $Succ(i) = \{j \in V \mid (i, j) \in E\}$ , i.e., the set of predecessors and of successors of  $i$ , respectively. For each task  $i$  let  $tl(i)$  denote the *top-level* of task  $i$ : if  $Pred(i) = \emptyset$ , then  $tl(i) = 0$ , otherwise

$$tl(i) = \max_{j \in Pred(i)} \{tl(j)\} .$$

Similarly, for each task  $i$  let  $bl(i)$  denote its bottom-level: if  $Succ(i) = \emptyset$ , then  $bl(i) = 0$ , otherwise

$$bl(i) = \max_{j \in Succ(i)} \{a_i + bl(j)\} .$$

Let  $pa(i, j)$  denote the length (as a sum of task weights) of the longest path from task  $i$  to task  $j$ , if such a path exists, otherwise let  $pa(i, j) = -\infty$ . The longest path length in  $G$  is then defined as  $d(G) = \max_{i, j \in V} \{pa(i, j)\}$ . Let us call  $d(G)$  the *failure-free makespan* of  $G$ . Because  $G$  is acyclic, we can compute  $d(G)$  in  $O(|V| + |E|)$  time as follows: add two zero-weight vertices  $v_1$  and  $v_2$  to  $G$ , where  $v_1$  represents a unique source task and  $v_2$  a unique sink task. Also add an edge from  $v_1$  to any entry task in  $G$  (a task without predecessor), and an edge from any exit task in  $G$  (a task without successor) to  $v_2$ . Then  $d(G) = pa(s_1, s_2)$  can be computed in time  $O(|V| + |E|)$  [42, Section 24.2].

We consider that tasks fail independently and task failure arrival times are exponentially distributed with Mean Time Between Failure (MTBF)  $1/\lambda$ . Therefore, the probability that task  $i$  fails during its first execution

attempt is  $1 - e^{-\lambda a_i}$ , in which case the task must be re-executed from scratch.

Our objective is to compute an approximation of the *expected makespan* of  $G$ , i.e., the longest path length in  $G$  taking into accounts that tasks can fail and must be re-executed. The failure-free makespan defined above is a clear lower bound on the expected makespan.

### IV. APPROXIMATING THE EXPECTED MAKESPAN

In this section we compute a *first-order approximation* of the expected makespan of a DAG  $G$ , which we denote as  $\mathcal{E}(G)$ . Our approximation relies on the fact that in practice  $\lambda$  is close to zero, which allows us to neglect  $O(\lambda^2)$  terms.

The probability that the first execution attempt of task  $i$  succeeds is:

$$p_i = e^{-\lambda a_i} = 1 - \lambda a_i + O(\lambda^2) .$$

The probability that the first execution attempt of task  $i$  fails but that its second execution attempt succeeds is :

$$(1 - e^{-\lambda a_i})e^{-\lambda a_i} = \lambda a_i + O(\lambda^2) = 1 - p_i + O(\lambda^2) .$$

Neglecting the  $O(\lambda^2)$  terms leads to the approximation that a task either takes time  $a_i$ , with probability  $1 - \lambda a_i$ , or time  $2a_i$ , with probability  $\lambda a_i$ . In other terms, our first-order approximation consists in assuming that a task never fails more than once. Hereafter when we say that a task fails, we mean that its first execution attempt fails and that its second execution attempt succeeds.

For any  $S \subset V$ , let  $P(S)$  denote the probability that all tasks in  $S$  fail and that no task in  $V \setminus S$  fails. Let also  $L(S)$  denote the length of the longest path in  $G$  when all tasks in  $S$  fail and no task in  $V \setminus S$  fails.  $\mathcal{E}(G)$  is thus defined as:

$$\mathcal{E}(G) = \sum_{S \subset V} P(S) \times L(S) .$$

We note that:

$$\begin{aligned} P(\emptyset) &= \prod_{i \in V} (1 - \lambda a_i + O(\lambda^2)) = 1 - \sum_{i \in V} \lambda a_i + O(\lambda^2) , \\ P(\{i\}) &= (\lambda a_i + O(\lambda^2)) \times \prod_{j \in V \setminus \{i\}} (1 - \lambda a_j + O(\lambda^2)) \\ &= \lambda a_i + O(\lambda^2) , \text{ and} \\ P(S) &= O(\lambda^2) \text{ if } |S| > 1 . \end{aligned}$$

Therefore:

$$\mathcal{E}(G) = P(\emptyset) \times L(\emptyset) + \sum_{i \in V} P(\{i\}) \times L(\{i\}) + O(\lambda^2) .$$



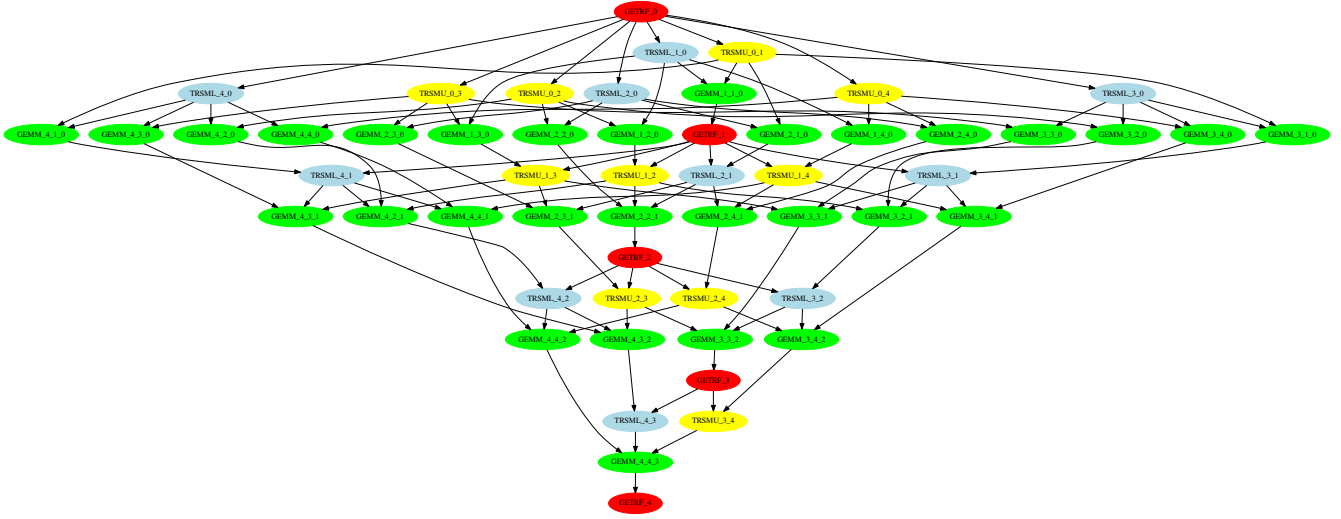


Figure 2: DAG of a LU factorization on a  $5 \times 5$  tiled matrix

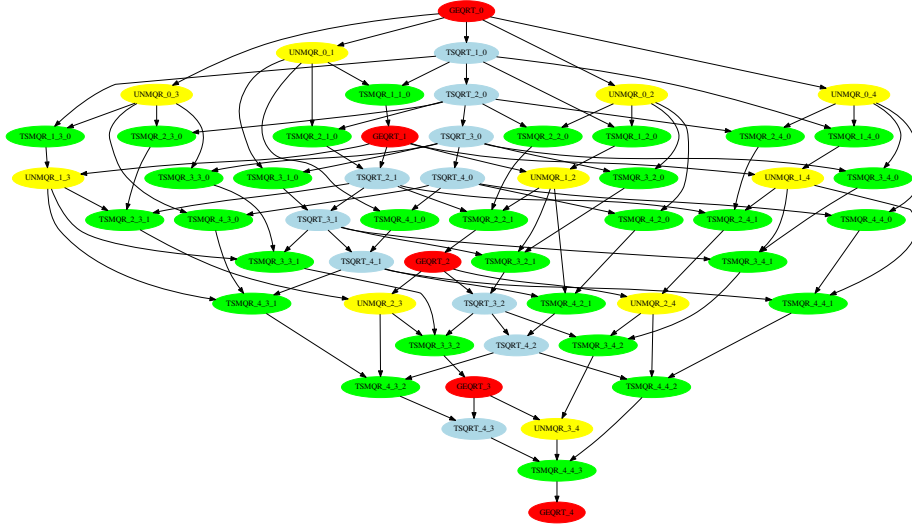


Figure 3: DAG of a QR factorization on a  $5 \times 5$  tiled matrix

To allow for consistent comparisons of results across different DAGs (with different number of tasks and different task weights), in our experiments we simply fix the probability that a task of average weight fails,  $p_{\text{fail}}$ , to 0.01, 0.001, or 0.0001. More precisely, for a given DAG  $G = (V, E)$  and a given  $p_{\text{fail}}$  value, we compute the average task weight as  $\bar{a} = \sum_{i \in V} a_i / |V|$  and pick the failure rate  $\lambda$  such that

$$p_{\text{fail}} = 1 - e^{-\lambda \bar{a}}.$$

While the above values of  $p_{\text{fail}}$  are chosen to evaluate the performance of our proposed approximation and to

compare it to competitors, we point out that these values are quite pessimistic. In other words, they lead to error rates much higher than those observed or expected on current and future large-scale computing platforms. The average execution time of a task in our experiments is  $\bar{a} = 0.15$  seconds, which leads to an error rate  $\lambda = 0.067$  with  $p_{\text{fail}} = 0.01$ . Equivalently, the MTBF is  $\mu = 1/\lambda = 14.9$  seconds. For a platform with 100,000 processors, this corresponds to an individual MTBF (per processor) of 17.27 days, quite an unrealistic value. The intermediate value of  $p_{\text{fail}}$ , namely 0.001, leads to an individual MTBF of 174 days. The smallest value of

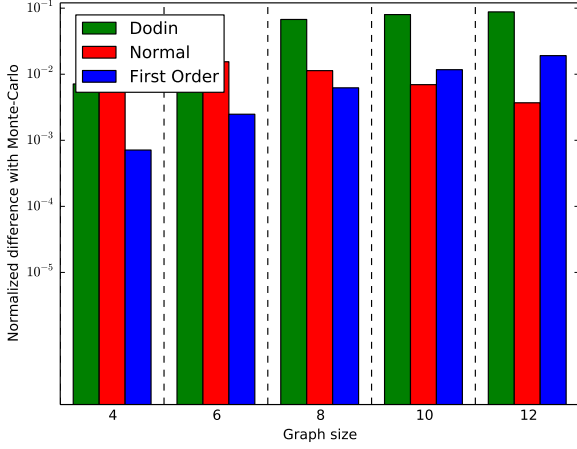


Figure 4: Cholesky,  $p_{\text{fail}} = 0.01$

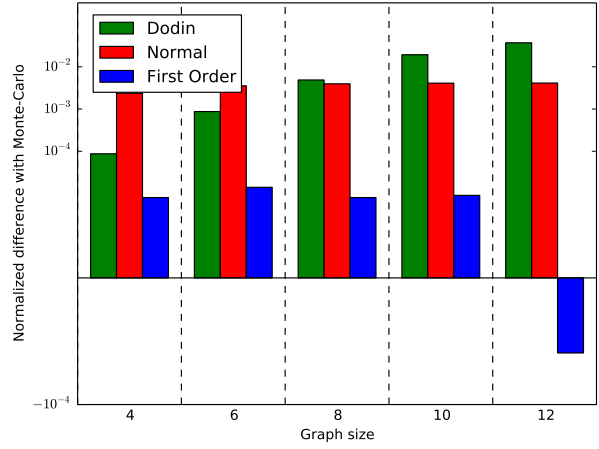


Figure 6: Cholesky,  $p_{\text{fail}} = 0.0001$

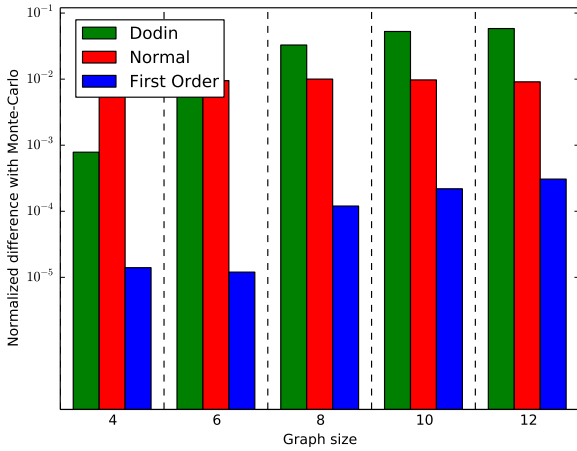


Figure 5: Cholesky,  $p_{\text{fail}} = 0.001$

$p_{\text{fail}}$ , namely 0.0001, leads to an individual MTBF of 4.7 years. As seen in the results in the next section, the lower  $p_{\text{fail}}$ , the lower the error incurred by our proposed approximation (and importantly much lower than that of its competitors). In other words, our selected relatively high  $p_{\text{fail}}$  values put our algorithm at a disadvantage with respect to its competitors.

#### D. Approximation Error Results

In this section we show relative error results, relative to expected makespans computed using the Monte Carlo method. Negative values denote an underestimation, while positive values denote an overestimation. All figures use a logarithmic scale on the vertical axis.

Figures 4, 5, and 6 show relative error vs. graph size for Cholesky graphs and for  $p_{\text{fail}} = 0.01, 0.001,$  and

0.0001, respectively. For  $p_{\text{fail}} = 0.01$ , we see that First Order leads to the lowest relative error for graphs sizes below 10, and that Normal leads to the lowest relative error for larger graphs. Dodin leads to the largest relative error, but for the smallest graph size. Considering the largest graph size, First Order has 1.9% relative error while Normal has 0.3% relative error and Dodin has 8.7% relative error. For  $p_{\text{fail}} = 0.001$ , First Order leads to dramatically lower error than its competitors (at least one order of magnitude lower). For instance, for the largest graph, First Order has 0.03% relative error, compared to 5.8% for Dodin and 0.9% for Normal. These results are even more striking for  $p_{\text{fail}} = 0.0001$ . In this case, still for the largest graph, First Order has -0.0006% relative error, compared to 3.7% for Dodin and 0.4% for Normal. In this case, First Order is an underestimation of the expected makespan.

Figures 7, 8, and 9 show similar results for LU DAGs. The overall message from these results is the same, meaning that First Order leads to drastically lower error than its competitors as  $p_{\text{fail}}$  decreases. Here again Dodin leads to the largest errors overall. For  $p_{\text{fail}} = 0.01$ , our approach leads to similar errors of the same order of magnitude as the errors for Normal (and in this case Normal is most often an underestimation while First Order is an overestimation).

Finally, Figures 10, 11, and 12 show results for QR DAGs and similar trends are observed. First Order drastically improves on Dodin and Normal for  $p_{\text{fail}} = 0.001$  and  $p_{\text{fail}} = 0.0001$ . For  $p_{\text{fail}} = 0.01$ , First Order leads to similar or lower error than Normal. Here again, Dodin leads to the highest errors across the board.

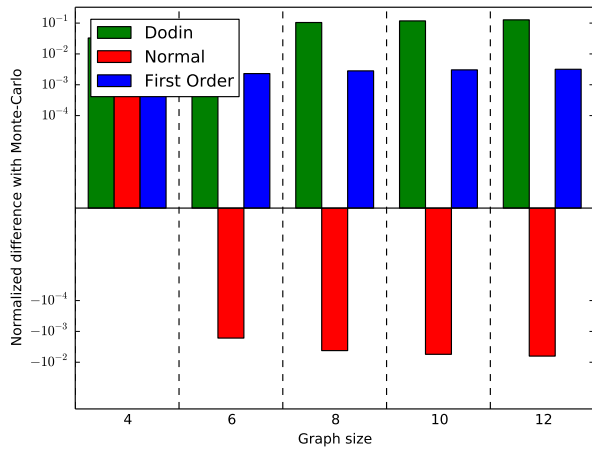


Figure 7: LU,  $p_{\text{fail}} = 0.01$

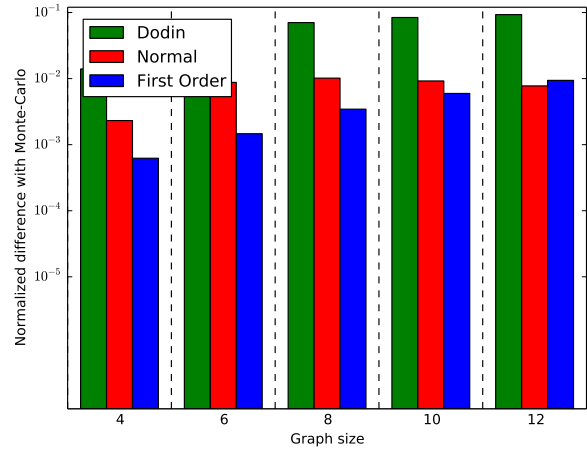


Figure 10: QR,  $p_{\text{fail}} = 0.01$

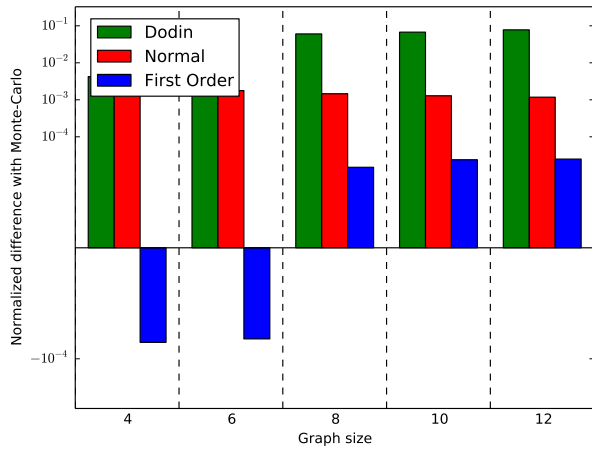


Figure 8: LU,  $p_{\text{fail}} = 0.001$

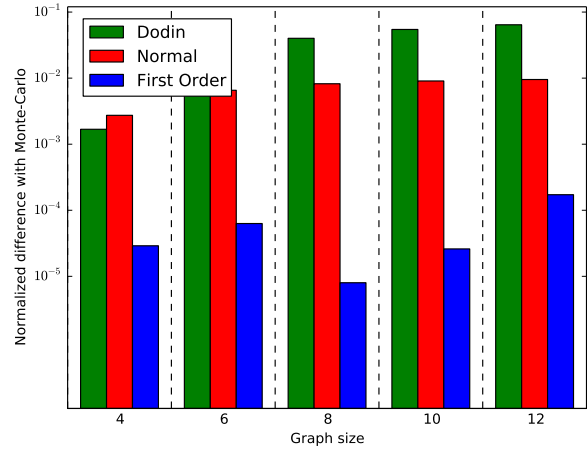


Figure 11: QR,  $p_{\text{fail}} = 0.001$

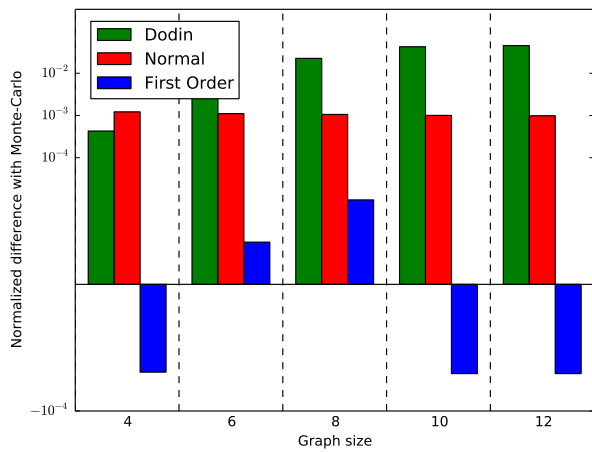


Figure 9: LU,  $p_{\text{fail}} = 0.0001$

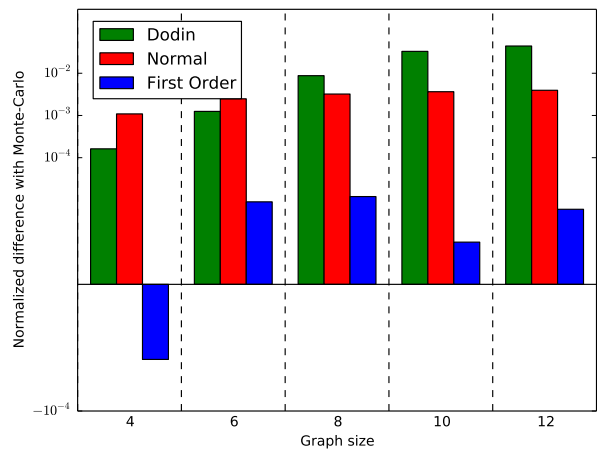


Figure 12: QR,  $p_{\text{fail}} = 0.0001$



### E. Scalability Results

To assess the scalability of the three approximation methods, we run experiments with  $k = 20$  for the LU DAG, i.e., 2,870 tasks, and  $p_{\text{fail}} = 0.0001$ . For this large graph, we ran the Monte Carlo for ten hours, so as to ensure the accuracy of the ground truth. Error (normalized difference with Monte Carlo) and execution time results are shown in Table I. We see that for this large DAG Dodin exhibits very large error. First Order is roughly 2 orders of magnitude more accurate than Normal. We also see that First Order can be computed within one second, while Normal requires about 20 minutes, i.e., about 3 orders of magnitude longer.

### F. Summary of Results

Our results show that our proposed approximation, First Order, is accurate as long as the probability of task failure is sufficiently low, and in particular much more accurate than its two competitors. When the task failure probability is high, its accuracy is comparable to or better than the Normal approximation. Across the board the Dodin approximation leads to high error. This is because the DAGs that we consider are far from being series-parallel. As a result, the series-parallel graph constructed by Dodin is a poor approximation of the original DAG, hence the large errors. Finally, not only is First Order more accurate, but also it is faster. Its execution on the large problem with 2,870 tasks requires negligible time, as opposed to several minutes for Dodin and Normal.

## VI. CONCLUSION

We have proposed an algorithm to compute a first-order approximation of the expected makespan of a Directed Acyclic Graph (DAG) of tasks in which tasks are subject to silent errors that may require task re-execution from scratch due to corrupted results, which we term a failure. Our approximation can be computed in  $O(V^2 + V.E)$  time for a DAG with  $V$  vertices and  $E$  edges, and it is exact in the first order and neglects second order  $O(\lambda^2)$  terms where  $\lambda$  is the exponential failure rate. This amounts to assuming that each task may need to be re-executed at most once. The problem of computing the expected makespan of a DAG of tasks with this assumption is actually a computationally difficult problem (#P-complete). As a result, techniques to approximate the expected makespan have been proposed in previous works [26], [24]. We have evaluated our proposed approximation and these previously proposed techniques for 3 classes of application graphs from the field of numerical linear algebra. These evaluations

quantify the approximation error with respect to a ground truth computed via a brute-force Monte Carlo method.

Our results show that our proposed approximation outperforms its competitors, by several orders of magnitude at low failure rates. In addition, it can be computed much more quickly than previously proposed approximations, which is crucially important for solving problems at scale. Overall, we have proposed a novel and improved approximation of the expected makespan of probabilistic DAGs in which tasks are subject to silent errors that mandate task re-executions.

Our general approach to obtain our first order approximation can be used to obtain a (more complicated but still tractable) second order approximation. While the improvement due to including the second order would be negligible for low failure rates, it may be significant for relatively high failure rates. A broader future direction is to adapt existing list scheduling scheduling algorithms, or to develop novel such algorithms, that rely on our proposed approximation to make scheduling decisions.

*Acknowledgments:* This work was supported in part by the French Research Agency (ANR) through the Rescue project. Yves Robert is with Institut Universitaire de France.

## REFERENCES

- [1] Y.-K. Kwok and I. Ahmad, "Benchmarking and comparison of the task graph scheduling algorithms," *J. Parallel Distrib. Comput.*, vol. 59, no. 3, pp. 381–422, 1999.
- [2] M. Wicczorek, R. Prodan, and T. Fahringer, "Scheduling of scientific workflows in the askalon grid environment," *SIGMOD Rec.*, vol. 34, no. 3, pp. 56–62, 2005.
- [3] E. Deelman, D. Gannon, M. Shields, and I. Taylor, "Workflows and e-science: An overview of workflow system features and capabilities," *Future Gener. Comput. Syst.*, vol. 25, no. 5, pp. 528–540, 2009.
- [4] P. Brucker, *Scheduling Algorithms*. Springer-Verlag, 2004.
- [5] Y. Robert and F. Vivien, Eds., *Introduction to Scheduling*. Chapman and Hall/CRC Press, 2009.
- [6] M. L. Pinedo, *Scheduling: Theory, Algorithms, and Systems*, 5th ed. Springer, 2016.
- [7] H. Topcuoglu, S. Hariri, and M. Y. Wu, "Performance-effective and low-complexity task scheduling for heterogeneous computing," *IEEE Trans. Parallel Distributed Systems*, vol. 13, no. 3, pp. 260–274, 2002.
- [8] F. Cappello, A. Geist, W. Gropp, S. Kale, B. Kramer, and M. Snir, "Toward exascale resilience: 2014 update," *Supercomputing frontiers and innovations*, vol. 1, no. 1, 2014.
- [9] A. Dixit and A. Wood, "The impact of new technology on soft error rates," in *IEEE International on Reliability Physics Symposium (IRPS)*, 2011, pp. 5B.4.1–5B.4.7.
- [10] D. Zhu, R. Melhem, and D. Mosse, "The effects of energy management on reliability in real-time embedded systems," in *Proceedings of the IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, 2004, pp. 35–40.

	Dodin	Normal	First Order
Normalized difference with Monte Carlo	-0.97	$954 \times 10^{-6}$	$7 \times 10^{-6}$
Execution time	around 2 minutes	around 20 minutes	less than 1 second

Table I: LU with  $k = 20$  and  $p_{\text{fail}} = 0.0001$ .

- [11] T. Héroult and Y. Robert, Eds., *Fault-Tolerance Techniques for High-Performance Computing*, ser. Computer Communications and Networks. Springer Verlag, 2015.
- [12] D. Fiala, F. Mueller, C. Engelmann, R. Riesen, K. Ferreira, and R. Brightwell, “Detection and correction of silent data corruption for large-scale high-performance computing,” in *Proc. of the ACM/IEEE SC Int. Conf.*, ser. SC ’12. IEEE Computer Society Press, 2012.
- [13] A. R. Benson, S. Schmit, and R. Schreiber, “Silent error detection in numerical time-stepping schemes.” *CoRR*, vol. abs/1312.2674, 2013.
- [14] Z. Chen, “Online-ABFT: An Online Algorithm Based Fault Tolerance Scheme for Soft Error Detection in Iterative Methods,” in *Proc. 18th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPOPP ’13. ACM, 2013, pp. 167–176.
- [15] P. Sao and R. Vuduc, “Self-stabilizing iterative solvers,” in *Proc. Scala ’13*. ACM, 2013.
- [16] G. Bosilca, R. Delmas, J. Dongarra, and J. Langou, “Algorithm-based fault tolerance applied to high performance computing,” *J. Parallel and Distributed Computing*, vol. 69, no. 4, pp. 410–416, 2009.
- [17] J. N. Hagstrom, “Computational complexity of pert problems,” *Networks*, vol. 18, no. 2, pp. 139–147, 1988.
- [18] L. G. Valiant, “The complexity of enumeration and reliability problems,” *SIAM J. Comput.*, vol. 8, no. 3, pp. 410–421, 1979.
- [19] J. S. Provan and M. O. Ball, “The complexity of counting cuts and of computing the probability that a graph is connected,” *SIAM J. Comp.*, vol. 12, no. 4, pp. 777–788, 1983.
- [20] H. L. Bodlaender and T. Wolle, “A note on the complexity of network reliability problems,” *IEEE Trans. Inf. Theory*, vol. 47, pp. 1971–1988, 2004.
- [21] R. H. Möhring, *Computational Discrete Mathematics: Advanced Lectures*. Springer, 2001, ch. Scheduling under Uncertainty: Bounding the Makespan Distribution, pp. 79–97.
- [22] M. Mitzenmacher and E. Upfal, *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press, 2005.
- [23] R. M. van Slyke, “Monte carlo methods and the pert problem,” *Operations Research*, vol. 11, no. 5, pp. 839–860, 1963.
- [24] L. C. Canon and E. Jeannot, “Correlation-aware heuristics for evaluating the distribution of the longest path length of a DAG with random weights,” *IEEE Trans. Parallel Distributed Systems*, 2016, available at <http://doi.ieeecomputersociety.org/10.1109/TPDS.2016.2528983>.
- [25] J. Valdes, R. E. Tarjan, and E. L. Lawler, “The recognition of series parallel digraphs,” in *Proc. 11th ACM Symp. on Theory of Computing*, ser. STOC ’79. ACM, 1979, pp. 1–12.
- [26] B. Dodin, “Bounding the project completion time distribution in PERT networks,” *Operations Research*, vol. 33, no. 4, pp. 862–881, 1985.
- [27] D. Sculli, “The completion time of PERT networks,” *The Journal of the Operational Research Society*, vol. 34, no. 2, pp. 155–158, 1983.
- [28] C. E. Clark, “The greatest of a finite set of random variables,” *Operations Research*, vol. 9, no. 2, pp. 145–162, 1961.
- [29] X. Ni, E. Meneses, N. Jain, and L. V. Kalé, “ACR: Automatic Checkpoint/Restart for Soft and Hard Error Protection,” in *Proc. SC’13*. ACM, 2013.
- [30] R. E. Lyons and W. Vanderkulk, “The use of triple-modular redundancy to improve computer reliability,” *IBM J. Res. Dev.*, vol. 6, no. 2, pp. 200–209, 1962.
- [31] J. Elliott, K. Kharbas, D. Fiala, F. Mueller, K. Ferreira, and C. Engelmann, “Combining partial redundancy and checkpointing for HPC,” in *Proc. ICDCS*, 2012.
- [32] K.-H. Huang and J. A. Abraham, “Algorithm-based fault tolerance for matrix operations,” *IEEE Trans. Comput.*, vol. 33, no. 6, pp. 518–528, 1984.
- [33] M. Shantharam, S. Srinivasmurthy, and P. Raghavan, “Fault tolerant preconditioned conjugate gradient for sparse linear system solution,” in *Proc. ICS ’12*. ACM, 2012.
- [34] M. Heroux and M. Hoemmen, “Fault-tolerant iterative methods via selective reliability,” Sandia National Laboratories, Research report SAND2011-3915 C, 2011.
- [35] G. Bronevetsky and B. de Supinski, “Soft error vulnerability of iterative linear algebra methods,” in *Proc. 22nd Int. Conf. on Supercomputing*, ser. ICS ’08. ACM, 2008, pp. 155–164.
- [36] E. Berrocal, L. Bautista-Gomez, S. Di, Z. Lan, and F. Cappello, “Lightweight silent data corruption detection based on runtime data analysis for HPC applications,” in *Proc. HPDC*, 2015.
- [37] L. Bautista Gomez and F. Cappello, “Detecting silent data corruption through data dynamic monitoring for scientific applications,” *SIGPLAN Notices*, vol. 49, no. 8, pp. 381–382, 2014.
- [38] —, “Detecting and correcting data corruption in stencil applications through multivariate interpolation,” in *Proc. 1st Int. Workshop on Fault Tolerant Systems (FTS)*, 2015.
- [39] B. Zhao, H. Aydin, and D. Zhu, “Reliability-aware dynamic voltage scaling for energy-constrained real-time embedded systems,” in *Proceedings of the IEEE International Conference on Computer Design (ICCD)*, 2008, pp. 633–639.
- [40] G. Aupy, A. Benoit, and Y. Robert, “Energy-aware scheduling under reliability and makespan constraints,” in *Proceedings of the International Conference on High Performance Computing (HiPC)*, 2012, pp. 1–10.
- [41] A. Das, A. Kumar, B. Veeravalli, C. Bolchini, and A. Miele, “Combined DVFS and mapping exploration for lifetime and soft-error susceptibility improvement in MPSoCs,” in *Proceedings of the Conference on Design, Automation & Test in Europe (DATE)*, 2014, pp. 61:1–61:6.
- [42] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*, 3rd ed. The MIT Press, 2009.
- [43] J. Choi, J. Dongarra, S. Ostrouchov, A. Petitet, D. Walker, and R. C. Whaley, “The design and implementation of the ScaLAPACK LU, QR, and Cholesky factorization routines,” *Scientific Programming*, vol. 5, pp. 173–184, 1996.
- [44] C. Augonnet, S. Thibault, R. Namyst, and P.-A. Wacrenier, “StarPU: A unified platform for task scheduling on heterogeneous multicore architectures,” *Concurr. Comput. : Pract. Exper.*, vol. 23, no. 2, pp. 187–198, 2011.