

Asynchronous and Exact Forward Recovery for Detected Errors in Iterative Solvers

Luc Jaulmes, Marc Casas, Miquel Moretó, Eduard Ayguadé, Jesús Labarta,
and Mateo Valero, *Fellow, IEEE*

Abstract—Current trends and projections show that faults in computer systems become increasingly common. Such errors may be detected, and possibly corrected transparently, e.g. by Error Correcting Codes (ECC). For a program to be fault-tolerant, it needs to also handle the Errors that are Detected and Uncorrected (DUE), such as an ECC encountering too many bit flips in a codeword. While correcting an error has an overhead in itself, it can also affect the progress of a program. The most generic technique, rolling back the program state to a previously taken checkpoint, sets back any progress done since then. Alternately, application specific techniques exist, such as restarting an iterative program with its latest iteration's values as initial guess.

We introduce a novel error correction technique for iterative linear solvers, designed to preserve both the progress made and the solver's future convergence by recovering the program's state exactly. Leveraging the asynchrony of task-based programming models, we mask our technique's overhead by overlapping error correction with the solver's normal workload. Our technique relies on analysing solvers to find redundancy in the form of relations between data. We are then able to restore discarded or corrupted data by recomputing or inverting the appropriate relations. We demonstrate that this approach allows to recover any part of three widely used Krylov subspace methods: CG, GMRES and BiCGStab, and their pre-conditioned versions.

We implement our technique for CG and recover lost data at the scale of a memory page, which is the granularity at which Operating Systems (OS) report memory errors on commodity hardware, and study the effect of varying the memory page size to address non-standard sizes and the possible use of huge pages in High Performance Computing (HPC). When compared to checkpointing and to the state-of-the-art algorithmic restart technique, on small (8 cores) to large scale (1024 cores), our methods show less overhead. A trade-off arises between our straightforward and asynchronous approaches, based on the rate at which faults happen. At the lowest considered rate and page size, overlapping recoveries decreases their average cost from 5.40% to 2.24% of the ideal faultless execution time. Our methods generally outperform the state-of-the-art even with increased overheads on big page sizes, and perform similarly on edge cases. These results also indicate that our techniques are increasingly efficient as the matrix size increases.

Index Terms—Resilience, Detected Uncorrected Errors, Exact Recovery, Krylov-Subspace Methods, HPC



1 INTRODUCTION

As High-Performance Computing (HPC) systems continuously scale up and chips shrink accordingly towards the exascale era [8], [9], the increased risk of faults in a computing system [41] causes reliability features to flourish in Operating Systems (OS) and processors [23], [28], [39]. Studies on actual software error rates indicate that the increase in faults due to decreasing size and increasing multiplicity of components is real, and that we can not rely solely on hardware to be resilient anymore [30], [35]. This tendency is expected to be aggravated by the reduction of voltages required for exascale systems [15], [26].

Memory cells are among the most vulnerable hardware components [33]. They are typically protected by Error Correcting Codes (ECC) [34] implemented at the hardware level. The widespread "SECDED" ECC can for example detect and correct all single bit flips. It can also detect, but not correct, when two bits are flipped in the same codeword. This kind of errors are called Detected and Uncorrected Errors (DUE), which are reported in specific registers in most modern architectures [3], [24]. When such a register is set, the processor generates a machine check exception

that can be caught at the OS level. Different kinds of DUE can be caught and signalled to the software stack and thus assimilated to some data loss, provided the architectural state is safe.

The responsibility of reacting against DUE is handed to the software stack. Some straightforward approaches like cancelling the affected process or relocating a faulty memory page to another physical location may be effective against low fault rates, but they are insufficient against the predicted rates that processors will suffer in the future. Also, very aggressive resilience strategies like process triplication are completely impractical unless we face very high fault rates [18]. Therefore, intermediate solutions that recompute an approximation of the lost data [29] or that save the process state in a checkpoint with a certain frequency have been extensively used [12], [32], [37]. However, most of these solutions involve backward recoveries, discarding useful computations, and thus incur significant slowdowns.

The application itself may be able to handle the error and terminate cleanly [5] or perform some sort of recovery procedure relying on Algorithmic-Based Fault Tolerance (ABFT), which has been extensively applied to MPI programs [10], [17], [29], as well as shared memory programming models [38], [40]. Algorithmic approaches have demonstrated to be more efficient than backward recoveries like checkpointing-rollback. However, most ABFT

- All authors are with the Barcelona Supercomputing Center (BSC)
E-mail: name.surname@bsc.es
- L. Jaulmes, M. Moretó, E. Ayguadé, J. Labarta and M. Valero are with Universitat Politècnica de Catalunya.

techniques so far suffer from two main drawbacks: being very application dependent, and still incurring significant overheads. In this paper we aim to reduce the impact of these two issues which have avoided the wide-spread usage of algorithmic resilience. The proposed ABFT methods to deal with DUE are based on very simple algebraic relations that do not require any kind of deep understanding of the algorithm and can be almost always derived for iterative methods. When a DUE is signalled, we always discard the whole memory page where affected data resides, as OSs do for bus and memory ECC errors, and recompute the discarded data using the algebraic relations that we derived. The cost of this recovery is proportional to the page size, and we evaluate the proposed ABFT methods for all possible sizes from 4K bytes up to 2M bytes. The overheads related to recovering data are reduced by overlapping them with algorithmic computations. Since the responsibility of such overlapping is left to the runtime system, we do not significantly increase the programming burden.

This paper proposes an integrated resilience approach, where the error detection is performed by hardware mechanisms that report DUE to the OS, which identifies lost data at a memory page level and triggers a signal caught by the application. We use the OmpSs task-based data-flow programming model [16], in which serial code is split into several pieces, called tasks, that are dynamically scheduled according to data dependencies explicitly expressed by the programmer. We combine the OmpSs annotations with MPI to scale our implementation up to over thousand cores. We demonstrate the feasibility of our approach by applying it to relevant iterative methods of the Krylov subspace family: the Conjugate Gradient (CG), Bi-Conjugate Gradient Stabilized (BiCGStab), and Generalized Minimal RESidual (GMRES) [4], and implementing it for CG. The main contributions of this paper are:

- A general resilience solution for DUE based on straightforward algorithmic recoveries. With the lowest error injection rate considered, corresponding to one expected error per baseline execution time, the overhead of this technique is 5.40%, whereas that of the checkpointing-rollback technique is close to 55%.
- An asynchronous and programmer transparent variant of our recovery implementation that reduces the overhead down to 2.24% under the lowest error rate, and that offers a trade-off between low overhead and convergence rate for higher error rates.
- A mathematical proof showing that Langou et al.'s Lossy Approach [29] is the best ABFT recovery strategy of all the restart techniques in the literature.
- An exhaustive and comprehensive evaluation, using real world matrices, of our method against a more sophisticated algorithm-specific restart method, derived from the Lossy Approach, and checkpointing-rollback mechanisms. We consider different parallel scenarios from 8 up to 1024 cores and we show that our methods always improve the performance of the above mentioned state-of-the-art methods.
- This paper further extends the original conference manuscript [25] with an in-depth study of the effect of page sizes, from 4KB up to 2MB, on the overheads

of the techniques. Our algorithmic methods outperform the state-of-the-art on average up to 512KB page sizes. For bigger sizes, our methods still perform better for the bigger matrices of the test set, and perform similarly to the Lossy Restart method on small matrices where a full vector can fit inside a single memory page.

The rest of this paper is organized as follows: Section 2 explains how to recover from hardware detected memory errors by using inherent redundancy, while Section 3 shows how to use this redundancy to make Krylov subspace methods resilient, as well as implementation details of this methodology for CG. Section 4 introduces the methods with which we compare our recoveries, and the next sections provide the experimental setup in Section 5, numerical validation in Section 6, and the evolution of these results with page sizes in Section 7. Section 8 examines related work and, finally, Section 9 provides our concluding remarks.

2 EXACT INTERPOLATION RECOVERY

2.1 Error Detection and Reporting

Due to the advent of faults, many processors have specific registers dedicated to signalling errors to the OS layer. On modern x86 and AMD64 architectures for example, a memory controller discovering data that is incoherent with the ECC, while accessing or periodically scrubbing it, reports it in a specific register [3], [24]. For memory pages, when the corrected errors exceed a threshold, the OS transparently relocates the page at another physical location. When a DUE is reported, the OS kills the affected process. This feature is known as memory page retirement on Solaris, and soft or hard page offlining in Linux kernels [28], [39].

In practice, application termination after a page failure is done by a SIGBUS signal. This signal can be caught and also specifies the failing memory addresses. By catching the signal and requesting a new hardware page at the same virtual address, the program can continue executing without further errors. Thus, to be resilient against memory DUE, an HPC application simply has to be able to replace lost data.

We classify all data as either constant if it does not change during execution time (matrix, preconditioners, right-hand side), or dynamic if it may be modified. Constant data is assumed to be saved to a reliable backing store, from which it is reloaded when errors are detected, similarly to other work using memory-page level fault models [7]. While there typically is not enough information available in the solver to recover constant data, there may likely be where the matrix was generated or read. Alternately, such data could be protected by software ECC at low cost, by exploiting the fact that this second ECC tier only needs to correct, and not detect errors [42]. Thus, only dynamic data needs to be made recoverable.

2.2 Extracting Redundancies of Linear Solvers

Linear iterative solvers perform operations like matrix-vector multiplications $q = Ap$, linear combinations $u = \alpha v + \beta w$, and combinations of the above, e.g. the very common residual $g = b - Ax$, where A is a matrix while q, p, u, w, g, b and x are vectors and α and β are scalars. In many cases the left and right hand side of these operations coexist during the whole execution of the solver.

TABLE 1

Block recoveries for operations $q = Ap$, $v = \alpha v + \beta w$ and $g = b - Ax$

Block relation, recover left side	Inverted relation, recover right side
$q_i = \sum_{j=0}^{n-1} A_{ij}p_j$	$A_{ii}p_i = q_i - \sum_{j \neq i} A_{ij}p_j$
$u_i = \alpha v_i + \beta w_i$	$w_i = (u_i - \alpha v_i)/\beta$
$g_i = b_i - \sum_{j=0}^{n-1} A_{ij}x_j$	$A_{ii}x_i = b_i - g_i - \sum_{j \neq i} A_{ij}x_j$

In some cases, we know that such a relation between vectors holds true (minus round-off errors) without having to recompute them. For example, if we define $x' = x + \alpha p$ and $g' = g - \alpha q$ with the above notations, then $g' = b - Ax'$. Finally, similar relations can hold true by construction, without ever having been computed. By analysing an iterative solver, we find redundancies expressed in terms of explicit or implicit relations between data.

Trivially, if any vector g, q, u is lost or partially corrupted, it can be recovered by recomputing the relation involving that vector. Given the inverses of A, α, β and other potential operands of a relation, it would also be possible to recover a lost or corrupted p, v, w, t or x . However, solving $Ap = q$ for p or $Ax = b - g$ for x is as expensive as running the whole iterative method. The matrix A can not be inverted either in the general case, due to numerical and computational considerations. However, recoveries based on such redundancy relations are applicable if only a small portion of the data structures involved in a relation is lost. This is the case with our error model since modern hardware is able to report errors at memory page level. In order to operate at such fine grain level, the redundancies must be expressed in terms of relations between small blocks of data.

2.3 Block Decomposition

The relations exposed previously, decomposed in n blocks, are listed in Table 1. We use the normal block relation to recover the left-hand side (lhs) of a relation, and the inverse of this relation for the right-hand side (rhs). In the event this inverse relation relies on a diagonal block of the matrix, we need to use a solver to recompute the lost data. If we know that a diagonal block is non-singular, e.g. when A is Symmetric Positive Definite (SPD), we solve the inverse block relations with a direct solver. Otherwise we solve this relation in the sense of least squares for the full columns of the matrix corresponding to the lost memory page as input, similarly to what Agullo et al. do for restart methods [2].

The formula for x_i 's recovery, shown at the bottom of the right hand side of Table 1 has been used by Chen [11] to recover the iterate, in complement of implicit checkpointing methods. Our approach requires no checkpoints however, as we protect all vectors with interpolation methods – as detailed in Section 3.1 for CG, BiCGStab, GMRES, and in Section 3.2 for their preconditioned variants. Exploiting these relations for recovery is a novel idea, since all previous work on making Krylov-subspace solvers fault-tolerant relies on a fail-stop failure model in a distributed memory environment. The required information to use redundancy of linear relations is then not available since corresponding parts of different vectors are lost simultaneously.

Furthermore, the granularity of the blocks of lost data in our recoveries is very different from the one in the context of process failure, which allows different and faster recoveries. Indeed, our block decomposition is dictated by the underlying layers (hardware detection, OS, runtime)

that do the DUE reporting. This means the block size that we use as granularity for recovery is a memory page. For a typical off-the-shelf machine, memory pages are 4K bytes, though bigger pages can be used in HPC settings – provided architectural support, such as the 2MB “huge pages” on x86. We evaluate our technique for any power of 2 sized page from 4KB to 2MB, thus for data losses from 512 to 262,144 double precision floating-point values.

Any DUE in our data protected by relations can thus be rectified by applying a small amount of computations, at worst factorizing a diagonal block of a matrix if one is used by that relation. This is a forward recovery scheme, since we can continue executing the program with our interpolated replacement data and the data that is not affected by the error. When A 's diagonal block is non-singular or a linear relation is used, we can even guarantee the exact same data as was lost for all relations (up to rounding errors), thus guarantee the same convergence rate as when the algorithm is not subject to faults. These recovery operations are usually small compared to the total computations, since matrix dimensions reach up to more than a million rows for real-life problems, as available in the University of Florida sparse matrix collection [14].

2.4 Dealing with Multiple Errors

Our approach requires no assumptions on simultaneous errors. Indeed, our techniques can easily handle multiple errors (discovered simultaneously) in most situations. Errors can always be recovered if they affect different instances of blocked linear relations expressed in Table 1. However, if simultaneous errors impact a single relation we have two possible scenarios:

- 1) *Simultaneous errors in a single vector* are not a problem for our recovery strategy. This is trivial for vectors recovered from linear relations, and straightforward for submatrix relations [29]. For two failed blocks i and j , we can combine both block relations:

$$\begin{pmatrix} A_{ii} & A_{ij} \\ A_{ji} & A_{jj} \end{pmatrix} \begin{pmatrix} x_i \\ x_j \end{pmatrix} = \begin{pmatrix} b_i - g_i - \sum_{k \neq i,j} A_{ik}x_k \\ b_j - g_j - \sum_{k \neq i,j} A_{jk}x_k \end{pmatrix}$$
 This relation is extendable to any number of blocks, with an increasing submatrix size to be factorized.
- 2) *Simultaneous errors on related data*, e.g. both q_i and p_i for a given i in a $q = Ap$ relation. Assuming no other relationship allows to recover these data, we may fall back to a restart method, e.g. the Lossy Restart which is adapted from the Lossy Approach [29] to fit our error model (see Section 4.3).

In conclusion, our forward interpolation recovery relies on very simple redundancy relations that are easy to identify in any iterative method, and that can efficiently be used at memory page level. This recovery can deal with multiple errors, but that may imply a more expensive computation or, at worst, the usage of a restart method as fallback.

3 APPLYING RECOVERIES TO ITERATIVE SOLVERS

3.1 Making Redundancies Explicit

DUE are reported when a faulty operation is made or when trying to access data that is corrupted. Even for data corruption discovered by the OS while periodically

Listing 1. CG pseudo code with redundancy relations

1	$\epsilon_{old} \leftarrow +\infty$	
2	$g \leftarrow b - Ax$	
3	for t in $0..t_{max}$	
4	$\epsilon \leftarrow \ g\ ^2$	$g = b - Ax$
5	if $\epsilon < tol$: break	
6	$\beta \leftarrow \epsilon/\epsilon_{old}$	
7	$d \leftarrow \beta d + g$	$d = A^{-1}q \quad g = b - Ax$
8	$q \leftarrow Ad$	see 3.1.1
9	$\alpha \leftarrow \epsilon / \langle q, d \rangle$	$q = Ad \quad d = A^{-1}q$
10	$x \leftarrow x + \alpha d$	$d = A^{-1}q \quad x = A^{-1}(b - g)$
11	$g \leftarrow g - \alpha q$	$q = Ad \quad g = b - Ax$
12	$\epsilon_{old} \leftarrow \epsilon$	

scrubbing memory pages, no error is signalled until that data is accessed, in the hope the page will be freed.

So in order to make an iterative solver resilient with our technique, it is sufficient to find for each operand of each operation done by the solver a relation

- that either allows to recover the operand, and then compute the result of the operation,
- or that allows to compute the result without this operand – that is, finding an alternate formulation.

Note that the main difference between this work and previous application-level recoveries for the same iterative solvers is the error model: since we do not consider complete failure of a node, we do not incur the loss of a part of every vector, which would render the relations we use here inapplicable. Over the next Sections we explain in detail how three commonly used iterative methods, CG, BiCGStab, and GMRES, can be protected using redundancy relations.

3.1.1 Conjugate Gradient

The pseudo-code for CG is given in Listing 1 [36], with the relations used for recovering each accessed data annotated on the right. Relations are written as whole-matrix relations for the sake of readability, but we use the memory page grained system described previously. Whenever possible, the relation that last produced data is used, which is not possible when data is updated in place. By construction, the algorithm conserves the relation $g = b - Ax$, and we can define an alternate way of computing q besides performing Ad from the update formula of $d \leftarrow \beta d + g$, which is $q \leftarrow \beta q + Ag$. When computing q , if a page of d is missing, q can be computed using this alternate formulation. When the whole matrix-vector multiplication is done, d can be recovered using $d = A^{-1}q$, in order to continue computations. However, it is impossible to use this relation when updating d . Let us consider the blocked formulation:

$$d_i = A_{ii}^{-1} \left(q_i - \sum_{j \neq i} A_{ij} d_j \right)$$

At this stage in the update of d_i , all pages $d_0..d_{i-1}$ are at iteration $t + 1$, and all pages $d_{i+1}..d_{n-1}$ are at iteration t .

We have two possibilities for recovery:

- 1) Postpone the recovery and compute q_i at iteration $t + 1$ using $\beta q + Ag$, then factorize A_{ii} to get d .
- 2) Perform *double buffering*, thus have two copies of a vector, either d or q , and use them alternately from one iteration to the next to remove in place updates.

The first option, though possibly more elegant, would imply taking rather considerable distances from the original algorithm. It is also arguable that the $\beta q + Ag$ operation might need to be protected, since a matrix-vector multiplication is the most computationally intensive and longest operation in a CG iteration. We thus opted for the latter option, and unrolled the loop to use two d vectors alternately, as illustrated in Listing 2. Code unchanged by this transformation is skipped. This solution adds redundancy to the method at the cost of some minimal memory overhead.

3.1.2 Bi-Conjugate Gradient Stabilized

BiCGStab is one of the generalizations of CG to matrices that are non-SPD. The pseudo-code for this method and the relations that may be used to make it resilient are presented in Listing 3, similarly to what has been done for CG. r is constant, along with the usual A and b . BiCGStab exhibits more redundancies than CG, and only an example set of relations that can be used is shown.

With $q = Ad$, $s = g - \alpha q$ and $t = As$, updating g can be rewritten $g \leftarrow g - \alpha Ad - \omega As$. Thus we have another way of computing g if for example q is faulty, but we also verified that the algorithm still conserves $g = b - Ax$.

Note that other assignments can also be expressed as slightly more complicated updates, we have e.g. $s \leftarrow s - \omega t - \alpha q$. The reverse also holds true, from the update of x we may get a direct relation such as $x = A^{-1}(b - s + \omega t)$

3.1.3 Generalized Minimal RESidual

The code for GMRES is available in Listing 4. Each iteration of GMRES consists of running the Arnoldi method - the part creating an orthogonal basis of vectors spanning $(g, Ag, \dots, A^{(m-1)}g)$ and an associated upper-Hessenberg matrix H - followed by a QR decomposition of this matrix H through Givens rotations. We may then increment the iterate by the solution y of $\min_y \|g - Hy\|$.

Protecting the biggest part of the data, which is the v^k vectors, is straightforward thanks to the Hessenberg matrix. At any time, we have at step t ,

$$l > 0 \text{ and } l < t \Rightarrow v^l = \frac{1}{h_{l,l-1}} \left(Av^{l-1} - \sum_{k=0}^{l-1} h_{k,l-1} v^k \right)$$

Thus the redundancy kept in the Hessenberg matrix' elements allows us to recover any Arnoldi vector under our error model.

Note that it is possible (and usual) to build the QR decomposition of the Hessenberg matrix H as the Arnoldi method goes, by computing the Givens rotation that corresponds to each new vector of the Arnoldi method. Q is thus computed as the set of Givens rotations, and $Q^T \|g\|_{2e_1}$ is also updated at every step. Thus we could use the relation $H = QR$ by keeping a copy of H even while building R :

- Givens rotations are easily deducible from H , thus Q and R are recoverable from H
- Givens rotations are easily invertible, since inverting a rotation means rotating by the opposite angle. Thus H is recoverable from Q and R .

Even though space is a limiting factor in GMRES, the H and R matrices are respectively upper Hessenberg and upper triangular of size $m(m + 1)$, thus much smaller than

Listing 2. CG with d double-buffered

```

1 for t in 0..t_max
2   ...
3   d1 ← βd2 + g
4   q ← Ad1
5   α ← ε / < q, d1 >
6   x ← x + αd1
7   ...
8   t++
9   d2 ← βd1 + g
10  q ← Ad2
11  α ← ε / < q, d2 >
12  x ← x + αd2
13  ...

```

Listing 3. BiCGStab pseudo code with redundancies listed

<pre> 1 g, r, d ← b - Ax 2 ρ ← < g, r > 3 for t in 0..t_max 4 q ← Ad 5 α ← ρ / < q, r > 6 s ← g - αq 7 t ← As 8 ω ← < t, s > / < t, t > 9 x ← x + αd + ωs 10 g ← s - ωt 11 check convergence 12 ρ_old ← ρ 13 ρ ← < g, r > 14 β ← ρ / ρ_old * α / ω 15 d ← g + β(d - ωq) </pre>	<p>d double-buffered</p> <pre> q = Ad g = b - Ax s = g - αq t = As x = A⁻¹(b - g) t = As </pre>	<pre> q = Ad s = A⁻¹t d = A⁻¹q s = A⁻¹t </pre>
	<pre> g = b - Ax q = Ad </pre>	<pre> d = A⁻¹q </pre>

Listing 4. GMRES pseudo code

```

1 for t in 0..t_max
2   g ← b - Ax
3   v0 ← g / ||g||2
4   for l in 0..m - 1
5     w ← Avl
6     for k in 0..l
7       hk,l ← < w, vk >
8       w ← w - hk,lvk
9     hl+1,l ← ||w||2
10    vl+1 ← w / hl+1,l
11    solve H = QR
12    y ← R-1QT||g||2e1
13    x ← x + ∑l=0m-1 ylvl
14    check convergence

```

the set of Arnoldi vectors of size mn (with $m \ll n$). Agullo et al. consider H to be stored (and solved) redundantly [2], which would then need no further protection. This also indicates that keeping the matrix H has a reasonable cost.

3.2 Preconditioned algorithms

The described recovery techniques can be straightforwardly applied to the same algorithms with a preconditioner. To preserve the generality of our approach, and to avoid preconditioners specifics, we consider a generic preconditioning operation “solve $Mu = v$ ”, M being the preconditioning matrix. To derive protected versions of the preconditioned algorithms we have to protect all the linear operations involving the preconditioned vectors. Protecting the execution of the preconditioner itself is beyond the scope of this paper, but a topic of complementary work, describing for example how to effectively protect multi-grid preconditioning [10].

To recover part of a preconditioned vector, there is no general way to avoid re-applying the preconditioner. Therefore, the prerequisite for the recovery to be cheap is the ability to perform a partial application of the preconditioner, that is, to apply the preconditioner to a small subset of v such that all lost data in u is recovered. If M is a block-diagonal matrix, solving $Mu = v$ only on the set of blocks that supersedes the lost data achieves this. If M is a fixed point method’s matrix, the sparse set of elements in v that contribute to the lost portion of u is sufficient. If M denotes a multigrid method, we consider the nodes of the coarsest grid that participate to producing lost data, then we only need the inputs that contribute to these nodes for recovery. In any case, re-running the preconditioner completely is a viable, though slow, forward recovery for u . Finally, a corrupted v after a “solve $Mu = v$ ” operation is always recoverable without using the equation $Mu = v$. This is an important point since M is not always explicitly formed.

This can be made explicit by looking at the preconditioned versions of CG, BiCGStab, and GMRES, which are shown in Listings 5, 6 and 7. We can easily observe that in both CG and BiCGStab the preconditioned vectors z, p and s always exist at the same time as their non-preconditioned counterparts, g, d and r , because the latter are still used in the solver. Thus we can always recover the preconditioned

vectors as discussed in the previous paragraph. All the relations protecting operations that involve z or g in CG, and p, s, d or r in BiCGStab are detailed next to the code of the preconditioned versions. For preconditioned GMRES, shown in Listing 7, the main redundancy relation from its non-preconditioned counterpart linking all the v^k is still valid. The only addition is the need for g to be conserved for the possible recovery of x .

3.3 Implementing Recovery with Asynchrony

CG and BiCGStab are harder to protect, as they require both redundancy relations and double buffering approaches to be fully protected, while GMRES just requires redundancies. For this reason, as well as because CG is a very popular method for solving SPD matrix equations in the HPC context, we select it to test our approach. We implement two versions of CG, one without a preconditioner and a second one using a block-Jacobi preconditioner. Any conclusion obtained from our experiments with CG can be trivially extended to the other two since they constitute a similar and simpler use-cases respectively, and to their preconditioned versions as explained in Section 3.2.

We start by presenting the implementation of CG for a shared-memory model, and extend it to a distributed memory systems in Section 3.4.

3.3.1 Conjugate Gradient’s Parallel Decomposition

The pseudo-code for CG is given in Listing 1, and its parallelization in tasks is done by strip-mining as shown in Figure 1a, with each set of tasks being named after the value or vector it outputs. Dependencies between tasks are generated from annotations to the sequential code, and represented by arrows on this graph. Tasks are then scheduled asynchronously by the runtime according to this data-flow. Some dependencies that do not affect the ordering or scheduling of tasks are not drawn for the sake of clarity.

Sets of tasks depicted in white represent operations that are strip-mined into as many parallel tasks as available threads. Blue tasks (after converging arrows) depend on all the previous tasks (because of a reduction operation) and represent a single task producing a scalar value. They are thus de facto synchronization points. The lattice-like arrows describe the fact that each following task depends on each

Listing 5. Preconditioned CG

```

1  $\epsilon_{old} \leftarrow +\infty$ 
2  $g \leftarrow b - Ax$ 
3 for  $t$  in  $0..t_{max}$ 
4   solve  $Mz = g$ 
5    $\rho \leftarrow \langle z, g \rangle$ 
6    $\beta \leftarrow \rho / \rho_{old}$ 
7    $d \leftarrow \beta d + z$ 
8    $q \leftarrow Ad$ 
9    $\alpha \leftarrow \epsilon / \langle q, d \rangle$ 
10   $x \leftarrow x + \alpha d$ 
11   $g \leftarrow g - \alpha q$ 
12   $\rho_{old} \leftarrow \rho$ 

```

Listing 6. Preconditioned BiCGStab

```

1  $g, r, d \leftarrow b - Ax$ 
2  $\rho \leftarrow \langle g, r \rangle$ 
3 for  $t$  in  $0..t_{max}$ 
4   solve  $Mp = d$ 
5    $q \leftarrow Ap$ 
6    $\alpha \leftarrow \rho / \langle q, r \rangle$ 
7    $r \leftarrow g - \alpha q$ 
8   solve  $Ms = r$ 
9    $t \leftarrow As$ 
10   $\omega \leftarrow \langle t, r \rangle / \langle t, t \rangle$ 
11   $x \leftarrow x + \alpha p + \omega s$ 
12   $g \leftarrow r - \omega t$ 
13   $\rho_{old} \leftarrow \rho$ 
14   $\rho \leftarrow \langle g, r \rangle$ 
15   $\beta \leftarrow \rho / \rho_{old} * \alpha / \omega$ 
16   $d \leftarrow g + \beta(d - \omega q)$ 

```

Listing 7. Preconditioned GMRES

```

1 for  $t$  in  $0..t_{max}$ 
2    $g \leftarrow b - Ax$ 
3   solve  $Mz = g$ 
4    $v^0 \leftarrow z / \|z\|_2$ 
5   for  $l$  in  $0..m - 1$ 
6      $u \leftarrow Av^l$ 
7     solve  $Mw = u$ 
8     for  $k$  in  $0..l$ 
9        $h_{k,l} \leftarrow \langle w, v^k \rangle$ 
10     $w \leftarrow w - h_{k,l} v^k$ 
11     $h_{l+1,l} \leftarrow \|w\|_2$ 
12     $v^{l+1} \leftarrow w / h_{l+1,l}$ 
13    solve  $H = QR$ 
14     $y \leftarrow R^{-1} Q^T \|z\|_2 e_1$ 
15     $x \leftarrow x + \sum_{l=0}^{m-1} y_l v^l$ 

```

previous task, as the block-row matrix-vector multiplication takes a whole vector as input for each single block as output.

3.3.2 Packing Recovery Tasks out of the Critical Path

We divide those relations in blocks as described in Section 2.2 and maintain an atomic bitmask (e.g. an `int`) per block of failure granularity, thus per memory page. Each data vector and task output is represented by a bit in this mask. Thus, if a task T works on a page p of a vector, it can check whether (one of) its inputs(s) is corrupted, and if so skip the computation while marking the bitmask with the bit representing T 's output. We similarly skip computations that have inputs whose computation was skipped, due to their inputs being corrupted. This is necessary as to keep track of when errors happen and avoid overwriting data potentially needed for recovery, and works especially well with linear relations (which are the majority of considered relations). There is a memory overhead directly proportional to the size of the linear system n to store this information.

Skipping computations is critical for reductions, because a floating point accumulation can be irretrievably corrupted by adding (or multiplying by) `+/-inf` or `nan`. Through a thread-private `sig_atomic_t` variable, each task is made aware of interruptions, and only contributes a page-level accumulation to the task-level one when no errors were reported. For this page based division to be valid, the reduction needs to be associative, which is always guaranteed since it is already required for the strip-mining into tasks.

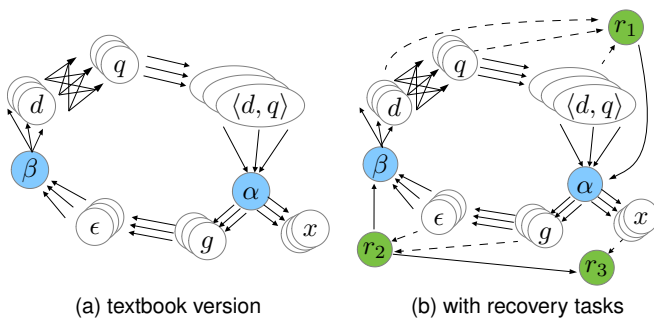


Fig. 1. Task decomposition of CG. Circles represent tasks producing the data inscribed on them, with white sets of tasks for strip-mined operations and blue for tasks producing a scalar, and full arrows representing data dependencies. Green tasks implement recoveries of the tasks linked to them with dashed arrows.

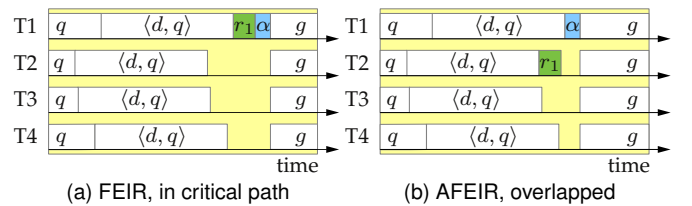


Fig. 2. Traces illustrating scheduling of recovery tasks

While errors are not corrected, the skipping of computations that depend on not-produced data propagates through the different tasks. When reaching a scalar task, skipping dependent computations would mean stop progressing completely. Thus we have to recover errors before the said scalar tasks. The graph in Figure 1b shows the modified cycle with the green tasks where recoveries take place: replacing lost data and recomputing skipped computations.

Each recovery task recovers the inputs and outputs of normal tasks which point to it with dashed lines. Recovery tasks are always added to the execution flow of the program and check the global variables for signalled errors. If none occurred, the recovery tasks do nothing.

The more conservative approach allows the recovery tasks to execute in the critical path, that is, waiting for all computations that do not need lost data to finish and only then run the recovery, as illustrated in Figure 2a. This option makes no compromise on the coverage of faults, since all the tasks (thus potential error discoveries) have finished executing when the recovery starts, as we will see with high error injection rates in Section 6.2. We call this technique Forward Exact Interpolation Recovery (FEIR).

Because dashed lines represent communication through atomic global bitmasks rather than dependencies, recovery tasks can execute concurrently to CG tasks, which do not touch the memory pages to be recovered, either skipping them or working on unrelated data. This allows to overlap computations and recoveries, thereby reducing overheads; however errors discovered between recovery tasks and the following scalar task are not recoverable. Thus we execute r_1 and r_2 asynchronously as late as the scheduler allows us to schedule them, which means concurrently with $\langle d, q \rangle$ and ϵ respectively, and with a lower priority as to start all reduction tasks first, see Figure 2b. Faults discovered after the recovery tasks start being executed can not be corrected.

This technique is the Asynchronous Forward Exact Interpolation Recovery (AFEIR).

The parallelization strategy for the preconditioned CG is exactly the same as for the non-preconditioned CG. The only added recovery technique is the partial “solve $Mu = v$ ”, as explained in Section 3.2, which is easy to perform since block-Jacobi is a blocked preconditioner.

3.4 Recovery on Distributed Memory Systems

The recovery methods described so far apply to shared memory models, in a single node or in a partitioned global address space for example. We list here the few modifications needed to extend our resilient methods to distributed memory programming models like MPI, which we will use to evaluate how the recovery techniques impact the scaling of the application in highly parallel scenarios.

More specifically, we use a hybrid implementation where the node level parallelism is leveraged using MPI and the intra-socket parallelism is exploited by the asynchronous task-based data-flow programming model, OmpSs. Our CG solver only requires the following additions:

- Global MPI reductions after the local reductions
- A task to exchange locally updated parts of p with neighbouring nodes depending on it.

This new exchange task takes place instead of the lattice-like dependencies between tasks d and q , while the MPI reductions occur during the α and β tasks (see Figure 1).

For our FEIR and AFEIR methods, we instantiate a second r_1 recovery task to be executed before our new exchange task, to avoid sending potentially failed and non-corrected data. Finally, MPI communications added inside the task recovering the x vector, r_3 , request and perform exchanges of parts of x when needed for recovery, since this vector is not exchanged at every iteration.

4 OTHER RECOVERY APPROACHES

4.1 Trivial Forward Recovery

The trivial forward recovery consists in simply keeping the program running, by allocating new (uninitialized) memory for corrupt or lost data. No other actions are taken. While an error in a part of the data that is not reused later would be masked, we lose all guarantees on convergence.

4.2 Rollback Recovery

Checkpointing is applied only to dynamic data in CG, to be fair in our comparison of methods, and to be consistent with the assumptions on which the DUE recovery relies. Each Processing Element (PE) periodically writes to its local disk the values of the iterate and search direction vectors it has at that given moment (x and p), which is the minimum to allow rolling back; frequency of checkpointing is expressed in terms of iterations of the solver.

There is no need to use a parallel file system, since we assume that the program will not crash (as we catch the errors). At rollback, each PE restores the vector portions that were saved to its local disk at the last checkpoint.

Checkpoints and rollbacks are global, that is, they involve all the PEs of a parallel run. In a distributed memory scenario, we perform a global MPI reduction once per iteration to decide whether a rollback is needed. This global communication is executed simultaneously with the α MPI reduction to avoid further synchronization overheads.

4.3 Lossy Restart

Langou et al. [29] present a forward recovery method for the fail-stop model of an MPI process, the Lossy Approach, applicable to all Krylov-subspace methods. To compensate for the loss of a part of the iterate x , a step of the block-Jacobi is used, which relies only on constant data and the remaining parts of x . This operation is similar to our recovery for the iterate, while discarding the residual in the block relation. After such an interpolation, a restart is necessary since the residual g is outdated and not easily deducible.

We adapt this Lossy Approach into a recovery for our error model, that we name the Lossy Restart:

- 1) If part of the iterate is lost we use the interpolation from the Lossy Approach. With i the failed block:

$$A_{ii}x_i = b_i - \sum_{j \neq i} A_{ij}x_j$$

- 2) We restart the method with either the intact or the newly interpolated iterate as initial guess.

Before comparing these methods, let us consider theoretical results presented on this interpolation, by noting x^* the solution of the system $b = Ax^*$, x the iterate, x^I the newly interpolated iterate, $e = x^* - x$ and $e^I = x^* - x^I$ the respective errors, and $g = b - Ax$ and $g^I = b - Ax^I$ the respective residuals. Langou et al. show the following:

Theorem 1. *The interpolation is contracting for a constant $c_i = (1 + \|A_{ii}^{-1}\| \sum_{j \neq i} \|A_{ij}\|)^{1/2}$, thus $\|e^I\| \leq c_i \|e\|$.*

This result has been improved for A SPD [2]:

Theorem 2. *With A symmetric positive definite, the interpolation diminishes the A -norm of the error: $\|e^I\|_A \leq \|e\|_A$.*

Both theorems grant the block-Jacobi’s fixed point property: if $x = x^*$, then $x^I = x^*$, since $e = 0$.

From here on, we will restrict ourselves to SPD matrices and show that the block-Jacobi step does not just give better replacement data, but the best possible – in the short run.

Theorem 3. *For A SPD, the interpolation minimizes the A -norm of the error $\|e^I\|_A$ over all possible values for x_i^I*

The proof of our theorem relies on the transformation of the error implied by the linear interpolation, $p_i^I : e \rightarrow e^I$ being a linear projection, orthogonal for the norm $\|\cdot\|_A$

Proof. By construction, the residual at x^I for the block i is $g_i^I = b_i - \sum_{j=0}^{n-1} A_{ij}x_j^I = 0$. Let us also notice that $g = b - Ax = A(x^* - x) = Ae$ and similarly $g^I = Ae^I$.

Now let us show that the kernel and image of p_i^I are orthogonal for A :

$$\begin{aligned} \forall e \in \mathbb{R}, \langle p_i^I(e), e - p_i^I(e) \rangle_A &= \langle e^I, e - e^I \rangle_A \\ &= \langle Ae^I, e - e^I \rangle \\ &= \sum_{j=0}^{n-1} \langle g_j^I, e_j - e_j^I \rangle \end{aligned}$$

This is always zero, since for $j = i$, $g_j^I = 0$ and for $j \neq i$, $x_j^I = x_j$ thus $e_j = e_j^I$. It then comes clearly that:

$$\|e\|_A = \|p_i^I(e)\|_A + \|(Id - p_i^I)(e)\|_A$$

where $p_i^I(e)$ depends solely on the e_j (thus x_j) with $j \neq i$. Hence the minimum of this norm for all possible x_i , or e_i , is reached in $p_i^I(e)$. \square

We can also deduce this from Theorem 2 and the fact that the unknown part of x is in the kernel of p_i^T . $\|e^T\|_A \leq \|e\|_A$ then holds for any x_i , hence the *min* relation of our theorem.

Restarting the solver, with a good or unmodified initial guess, still harms the superlinear convergence of CG, which relies on the fact that the sequence of iterates x minimizes at each iteration the norm $\|x^* - x\|_A$ on a sequence of increasing subspaces. However, this disturbance benefit methods who have a tendency to stagnate (such as GMRES).

All recoveries based on restarting are identical as long as the iterate is untouched, and trade in convergence properties for simplicity of recovery in the same way. It is to be expected that such methods would behave very similarly to the Lossy Restart, though always worse in the short run, hence it is the only restart method against which to compare.

5 EXPERIMENTAL SET-UP

We measure solving 9 matrices selected from the University of Florida sparse matrix collection [14]. They are well-conditioned matrices for CG selected among the biggest of each family of SPD matrices. We run experiments on an Intel® Xeon® E5-2670, with one thread on each of its 8 cores.

Our evaluation is done on two versions of CG, a non-preconditioned version to show the hardest case possible, and one using a block-Jacobi preconditioner. Due to the wide variety of preconditioners available for CG, it is impossible for us to evaluate every single one. We list in Section 3.2 the desirable properties of preconditioners for an efficient recovery. The block-Jacobi is simple to implement, and trivially applicable to a subset of a vector. We select it also because, if its block size coincides with the memory page size, the factorization of diagonal blocks for the recovery of single errors is already computed. Thus we will use diagonal blocks of 512 by 512 elements, which coincides with 4KB page sizes.

We compare the following methods: our Forward Exact Interpolation Recovery (FEIR) without asynchrony (recovery tasks in the critical path), our Asynchronous FEIR (AFEIR), the Lossy Restart, checkpointing-rollback to local disk, and trivial forward recovery. The optimal checkpointing rate is used whenever errors are injected, and no fallback is used for FEIR or AFEIR: simultaneous errors on related data are simply ignored (see Section 2.4).

5.1 Techniques Overheads

From here on, the “ideal” CG will refer to our version of CG with no resilience mechanisms nor error injections.

We present in Table 2 the harmonic means of overheads for all methods in absence of faults, compared to the ideal CG. The Lossy Restart and trivial techniques have no overhead when no errors are injected, since catching the error, replacing memory pages and ordering a restart is done in a signal handler which is never called. To give a sense of checkpointing cost we arbitrarily consider checkpointing periods of 1000 and 200 CG iterations. The corresponding overhead raises from 17.62% to 46.20% as the checkpointing frequency increases, which constitutes a significant cost.

The overheads associated to the AFEIR and FEIR techniques are much smaller since they are associated to activities like task creation or scheduling, that are much cheaper than writing data to disk. The asynchronous nature of

TABLE 2
Resilience methods’ overheads, no errors

	Lossy	trivial	AFEIR	FEIR	ckpt 1K	ckpt 200
overhead	0.00%	0.00%	0.23%	2.73%	17.62%	46.20%

TABLE 3
Increase of time spent per state for FEIR methods

	imbalance	runtime	useful
AFEIR	4.30%	8.11%	1.90%
FEIR	25.06%	7.84%	2.78%

the AFEIR technique allows to compensate much of the overhead incurred by the FEIR technique. We can see in Table 3 a detailed breakdown of what is involved in the overheads of the FEIR and AFEIR methods, expressed as the increase of the proportion of time spent in each state while the solver is running: either idle, thus suffering load imbalance, or performing runtime work, such as creating and scheduling tasks, or finally executing tasks, thus doing computations for the solver. Executing the recovery tasks in the critical path obviously increases the load imbalance.

Most of the runtime overhead of FEIR and AFEIR techniques could be removed if application-level resilience were supported by the runtime, instantiating recovery tasks only when DUE are signalled.

5.2 Error Injection

We consider the most common DUE to test the considered recovery techniques: the corruption of a memory page. However, this is generalizable to more types of errors, since a DUE very often ends up being a data corruption. DUE can also bring changes in the control flow of the programs, which typically ends up with a data corruption or an execution failure. Our model covers the first scenario, while the second lies beyond the scope of this paper.

Errors are injected from a separate thread at times defined by an exponential distribution parametrized by the Mean Time Between Errors (MTBE). To account for the wide range of convergence times across different matrices (from 1 to 100 seconds), we normalize the MTBE to the ideal convergence time. Affected memory pages are selected at random with uniform distribution.

To simulate errors we use the `mprotect` system call available in Linux kernels to change the authorizations of the targeted memory page. This is more practical than triggering a real hardware retiring of a memory page, and behaves identically: the program receives a signal at the time of access to the memory page. We recover in the same way as we would from a real error: in a signal handler, we request a new memory page at the same virtual address through means of the `mmap` system command. All the recoveries operate exactly in the same way as they would if a real DUE took place. For the solver, there is no difference between real hardware DUE and our error injection mechanism.

Errors are injected in the memory pages of the Krylov vectors, which we cover with our resilience techniques, and not in the pages that contain constant data, program instructions, scalar or control values of the algorithm, or our bitmasks used for resilience. The amount of non-constant data that is not covered by our error injections is very small compared to the targeted memory space, and relatively constant across the different resilience techniques.

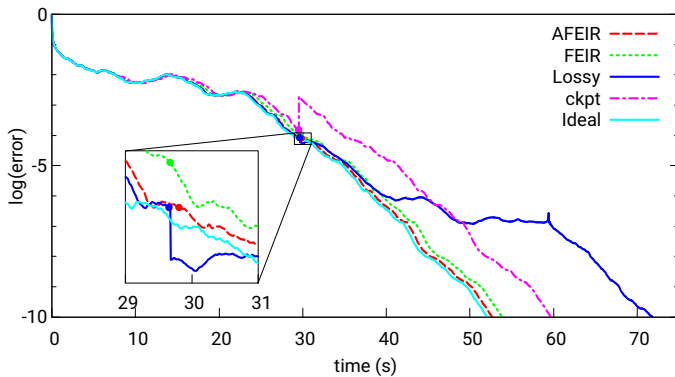


Fig. 3. CG convergence for different resilience methods with matrix thermal2 and same single error injection in the iterate x around $t = 30$ s

6 EVALUATION

6.1 Convergence

Figure 3 illustrates the convergence of CG for a sample scenario consisting of a single error injection. The x-axis represents the time and the y-axis shows the execution progress in terms of the logarithm of the residual norm defined as $\|Ax - b\|/\|b\|$, updated at each iteration. The ideal CG is represented by the cyan line; all the other experiments have a single error injected 30 seconds after the beginning of the execution at a certain memory page that contains a portion of the iterate x . Before the error is injected, each resilience method pays its typical overhead in absence of faults. The purple line corresponding to the checkpointing mechanism is the one with more overhead, which is consistent with the analysis presented in Section 5.1, as we use a checkpointing frequency of 1000 iterations. At the time of the error, it already incurs a 9.12% slowdown. Once the error is injected, the checkpointing mechanism rolls back a certain number of iterations and resumes progress from there. The Lossy Restart, represented by the blue line, has an immediate reduction in the error thanks to its block-Jacobi step interpolation, but converges slower afterwards because restarting harms CG’s superlinear convergence. The FEIR and AFEIR methods recover the lost data by using an exact interpolation and keep progressing. The overhead paid by the AFEIR technique is significantly smaller than the one paid by FEIR, since asynchrony allows most of the recovery work to be overlapped with other computations.

6.2 Shared-Memory Performance

Figure 4 shows an exhaustive evaluation of the performance slowdown associated to the 5 resilience mechanisms listed in Section 5: trivial, checkpointing-rollback, Lossy Restart, FEIR and AFEIR. The checkpointing rate is computed for each experiment to minimize execution time, taking into account the time to write and read checkpoints, the MTBE, and no downtime [6]. We consider the same 9 input matrices as for the overhead measures, and 6 error injection scenarios per matrix and method, which means that we provide an evaluation of 270 different experiments. Each experiment has been run over 50 times and Figure 4 reports their harmonic mean, and standard deviation as error bars. In each repetition, the errors have been injected randomly at different times and memory pages. On the x-axis of Figure 4 we display the name of the considered matrices and, for each matrix, the error injection frequency

normalized to the ideal CG’s convergence time τ for that matrix. A value n means an error frequency of $\frac{n}{\tau}$, thus an MTBE of $\frac{\tau}{n}$. In other words, n is the expected number of errors injected during the ideal convergence time τ . The y-axis is displayed in logarithmic scale and shows the measured performance slowdown in percentage for each experiment, with respect to the ideal CG. A slowdown close to 0 means the resilient CG converges at a speed close to that of the ideal one, whereas a bigger slowdown means its convergence is slower. The convergence threshold is 10^{-10} .

We have run the exact same 270 experiments with the block-Jacobi Preconditioned CG (PCG), however due to space limitations, we only report the mean of those results, displayed at the right hand side of Figure 4.

The trivial method reacts badly against few errors and its convergence times diverge extremely fast, with overheads over 200% with a normalized error frequency of 5 only. For PCG, the overheads of the trivial recovery reach 50% with a normalized frequency of 2 and become larger than 700% for frequencies of 10 or more. The checkpointing scheme reacts better than the trivial method, with substantial overheads that tend to increase slower, ranging on average from 55% to 433% for CG and from 60% to 752% for PCG. These convergence times are close to the expected values from the checkpointing frequency computation [6]. The Lossy Restart behaves better on average than the trivial method and the checkpointing schemes. Regarding CG, it has an overhead of 8.4% with one expected error per ideal convergence time, reaching up to 87% and 170% against 20 and 50 times higher frequencies respectively, whereas for PCG these overheads are 12.80% and 500% for normalized frequencies of 1 and 20. This better behaviour of the lossy mechanisms with respect to checkpointing and trivial techniques is already reported in the literature [2], and the fact that our experimental framework has reproduced known results demonstrates its accuracy and reliability. This is matrix specific however, as non-ABFT methods are clearly outperformed for `af_shell8` or `cfid2` but have overheads similar to Lossy for `thermal2`.

The most important fact of our evaluation is that methods FEIR and AFEIR behave much better than the current state-of-the-art resilience techniques for iterative solvers. When applied to CG, FEIR has an overhead of 5.37% and 29.68% under normalized frequencies of 1 and 50, whereas AFEIR has overheads of 3.59% and 50.47% respectively for the same error rates. On PCG, FEIR and AFEIR have an overhead of 5.36% and 2.72% with the smallest frequency, and reach 40.55% and 48.55% with the largest.

While FEIR has a roughly constant overhead on most matrices, the impact of recoveries in the critical path can be seen where execution times are the shortest, such as for `Dubcova3` and especially `qa8fm` and `thermomech`. As recoveries run on a per iteration basis, the error rates per iteration determine the chances of encountering errors on related data and during recovery tasks. Under injection frequencies of 20 and 50, both `qa8fm` and `thermomech` experience over 0.2 and 0.6 errors per iteration, which significantly pulls up the mean overheads. Such extreme cases could be dealt with by using more recoveries per iteration, or a fallback method for unrecoverable errors. These matrices also show that Lossy Restart is most efficient on fast converging problems.

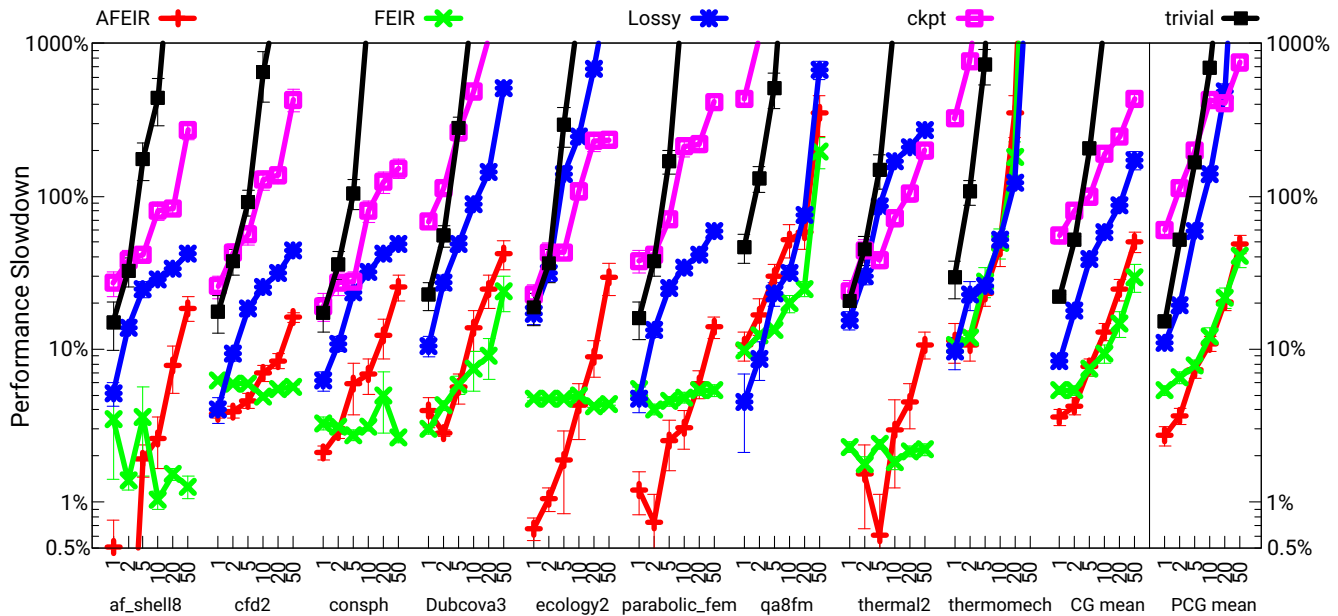


Fig. 4. Comparison of the execution time for resilience methods and matrices, varying error injection rates

AFEIR is slower than FEIR for high error injection rates, because errors happening between the end of a recovery task and its following scalar task are unrecoverable. That is the time between the end of r_1 or r_2 and the beginning of α and β respectively, as illustrated by Figure 2b. With very high error injection rates, the probability of an error happening during these time windows may cause the contribution of a memory page to $\langle d, q \rangle$ or ϵ (see Figure 1b) to be ignored. Depending on the matrix and the actual data lost, this might have a significant impact, as matrix ecology2’s behaviour shows. The FEIR method is not at risk of discovering an error after a recovery task ran, because these tasks start after all computations are done. However, both methods are still vulnerable during the recovery’s execution. There is thus a trade-off between the low overheads of AFEIR at frequencies of 10 and less, and a more conservative approach, FEIR, which trades in some convergence speed for safer recoveries and is thus useful at higher error rates. The same trade-off applies to the PCG results for low error injection rates. The precomputed factorization of diagonal blocks reduces recovery time, thus a block-Jacobi preconditioner weakens this trade-off for high error injection rates. It is to be expected that when using a preconditioner whose partial application is computationally hard (see Section 3.2), the average recovery time will increase and this trade-off will become stronger.

6.3 Scaling Results

In this section, an evaluation of the scalability of our recovery techniques is performed with a hybrid MPI + OmpSs implementation. Since the previously considered matrices are not well suited for large scale experiments, we solve Poisson’s equation in 3D using a 27 point stencil discretization, which is also used in the HPCG benchmark [20], with a system size of 512^3 unknowns. Experiments are run in the MareNostrum supercomputer, whose nodes contain two Intel Xeon CPU E5-2670 sockets. Each MPI rank is mapped to one 8-core socket, running 1 OmpSs thread per core. We

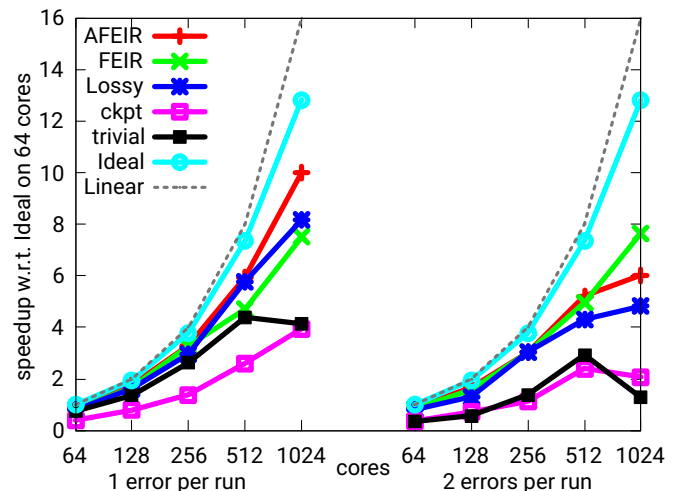


Fig. 5. Speedup of the MPI+OmpSs resilient CGs

consider runs on 8, 16, 32, 64 and 128 sockets (64 to 1024 cores), since we need 8 sockets to fit the matrix in memory.

We present in Figure 5 a complete evaluation in terms of speedup, injecting one and two errors per run. The speedups are computed taking the execution time of an ideal CG on the smallest possible core count, 64, as a reference. We display data concerning the FEIR, AFEIR, Lossy Restart, checkpointing and trivial techniques, and include the ideal CG’s and linear speedups for reference. Our MPI+OmpSs CG implementation achieves a parallel efficiency of 80.17% on 1024 cores in a faultless run, which highlights its quality in terms of parallel performance.

AFEIR and FEIR techniques clearly overcome the trivial and the checkpointing techniques, achieving speedups of 10.01 and 7.50 respectively when 1 error is injected and 6.03 and 7.65 against two errors on 1024 cores. The Lossy Restart achieves speedups of 8.17 and 4.82 respectively on 1024 cores. It is worth noting that only a few tens of iterations are required to achieve convergence for the 27-point stencil matrix, which causes any overhead to be important

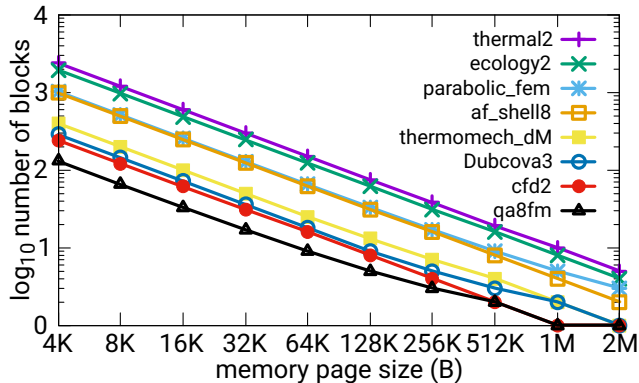


Fig. 6. Subdivision of CG vectors in blocks, with one block per page

compared to the ideal execution time, but also makes this matrix the ideal workload for a restart method. Even in this case, restarting is as costly if not more than our FEIR and AFEIR methods' overheads. Regarding checkpointing, writing vectors to disk already causes the checkpointing to perform significantly worse than our baseline, and when injecting errors its speedups stay below a third of the ideal CG, close to that of the trivial method.

7 ANALYSIS OF DATA LOSS GRANULARITY AND MEMORY PAGE SIZE

While the checkpointing scheme recovers full vectors, the algorithmic techniques recover or replace data at the memory page level, which is usually 4KB. However, some systems use different page sizes, and some HPC codes rely on huge pages (2MB) to improve performance, hence we present an in-depth study of the impact of varying page sizes on recoveries and convergence rates.

The impact of using huge pages on the ideal performance (without resilience methods nor error injections) was evaluated as well, however details are left out due to space constraints. Overall, the number of TLB misses decreased drastically with negligible impact on convergence time.

The number of pages, thus blocks of data erased by a single error, per vector (b, x, g, d, q , see Listing 1) are displayed in Figure 6. At the biggest page sizes, a single page may contain a vector entirely, thus full vectors can be erased by a single error, as is the case for qa8fm and cfd2 starting at 1MB page sizes, and Dubcova3 and thermomech_dM at a page size of 2MB. In this case, the Lossy recovery consists of solving the problem by factorizing directly the whole matrix until the resulting vector is not subject to faults.

7.1 Impact of page size depending on matrix size

To show the different impact of page size on different matrices, we display in Figure 7 the overheads for 3 different page sizes and 3 matrices: the smallest, qa8fm, the biggest, thermal2, and an intermediate one, af_shell8, which appear respectively at the bottom, middle and top of Figure 6. These matrices represent all possible behaviours that we observed. The plotted values are, as in previous Figures, average performance slowdown compared to the "ideal" baseline execution time, and this data is a subset of the data that will be presented in Figure 8. Each point corresponds to the harmonic mean of up to 70 runs divided by the ideal time, where each configuration consists of a recovery method, matrix, page size, and fault injection rate.

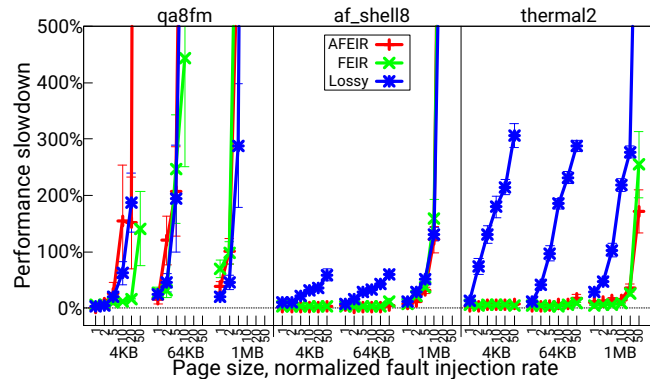


Fig. 7. Overheads of algorithmic techniques for 3 selected matrices and 3 page sizes, presented linearly

The overall tendency is the same as previously: FEIR and AFEIR perform significantly better than Lossy in most scenarios, with the notable exception of big page sizes for small matrices. However, looking at differences per matrix sizes, we note that all methods diverge for qa8fm under high fault injection rates, except FEIR with 4KB pages which incurs a 141% slowdown. For af_shell8, until 64KB pages the FEIR and AFEIR methods stay below 11% while Lossy peaks at 60%, and with 1MB pages the fault rates are similar until a fault injection rate of 5 faults per run, with quick divergence after that. The biggest matrix, thermal2, behaves similarly as FEIR and AFEIR have low slowdowns for page sizes of 4KB and 64KB while Lossy's slowdown seems roughly proportional to the fault injection rate, until 307%. With 1MB pages we see a similar divergence to previous matrices for Lossy, while the increase in slowdowns is slower for FEIR and AFEIR with 254% and 171% respectively.

This Figure allows to illustrate that the bigger the matrix, the better FEIR and AFEIR perform in comparison to Lossy. This is in part due to bigger matrices usually having worse convergence [36]. Such matrices will incur more slowdown from restarting than from exact recoveries, such as FEIR's and AFEIR's. Coincidentally, dividing a vector into more blocks due to a matrix being bigger also reduces the probability of having two faults affect related data, which adversely impacts FEIR and AFEIR as explained in Section 2.4.

7.2 Overall Page Size Evaluation

Figure 8 presents the average overheads of all algorithmic techniques (trivial, FEIR, AFEIR, Lossy) for CG, with different page sizes. We repeat for all page sizes ranging from 4KB to 2MB the full evaluation done in Figure 4 for each matrix, technique (except checkpointing), and fault injection rate, for a total of 2160 different configurations. We use the same methodology as explained in Section 5 and compute the harmonic means of slowdowns (convergence time divided by baseline execution time) for each configuration of up to 70 runs. We then report the harmonic means of all the matrices' slowdowns. That is, every point is the average of all the averaged runs of all the matrices, for that page size and fault injection rate. The first graph "4KB" coincides with the graph titled "CG mean" in Figure 4.

Overall, the trade-off between asynchronous and critical-path recoveries discussed in Section 6.2 is still applicable up to 64KB page sizes, above which FEIR and AFEIR show similar performances. The high cost of recoveries for 128KB

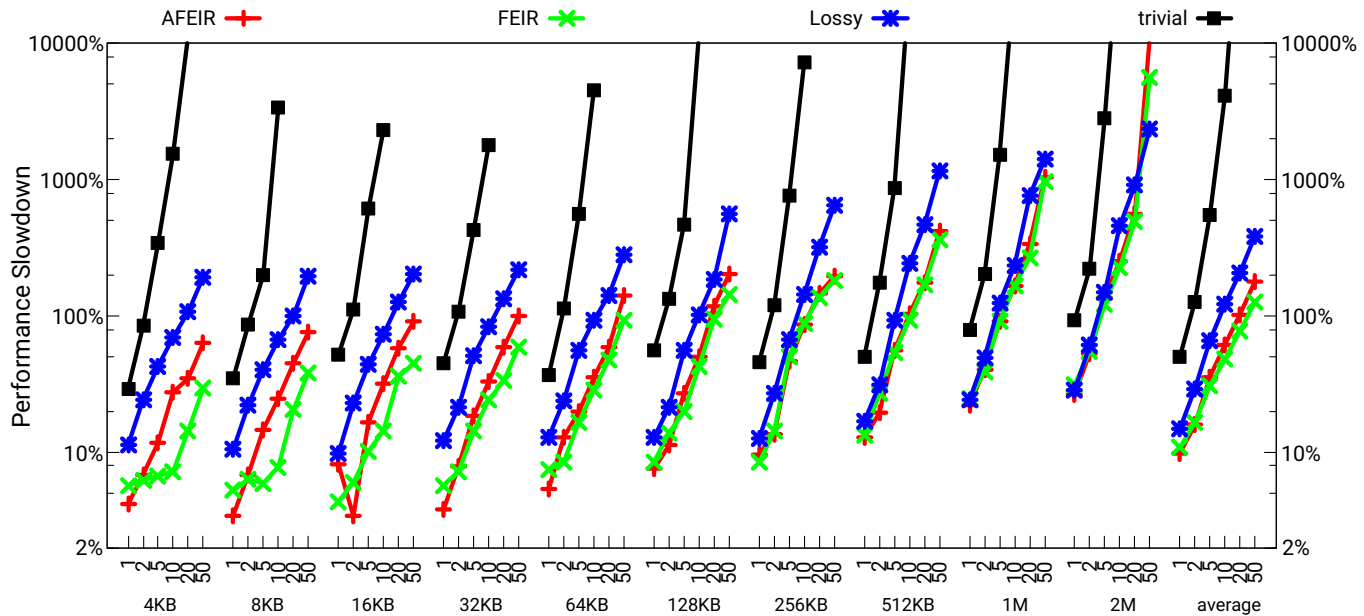


Fig. 8. Evolution of the overheads of algorithmic techniques for page sizes from 4KB to 2MB, varying error injection rates

page sizes and above makes the choice of overlapping recoveries largely irrelevant. Both methods outperform the Lossy Restart strategy significantly up to 512KB page sizes, above which all the convergence times are similar, and even to the advantage of Lossy for an error injection rate of 50 with 2MB page sizes.

At the same time, all forward recovery methods (FEIR, AFEIR and Lossy) still outperform by far the trivial method, even though it does not perform any computation. The trivial method often does not converge at all, as is the case with normalized error injection rates over 5 with page sizes of 128KB or more.

With increasing page sizes comes an increasing cost of computations during recoveries, and we can indeed see how increasing recovery costs drive convergence times up for all algorithmic techniques. The minimum slowdowns of AFEIR, FEIR and Lossy are all around 28% for a single expected fault per run at 2MB pages against 4.2%, 5.5% and 11.2% respectively with 4KB pages. This means all methods still outperform checkpointing, whose overhead for a single expected fault is 55% as presented in the previous Section. Maximum values increase, and for 1MB pages the slowdowns of FEIR, AFEIR and Lossy are 1072%, 1146% and 1516.59%. As in Section 6.2, these high numbers are mostly driven by the overheads of small matrices.

The different techniques are affected differently by the page size modification. For similar convergence rates, the Lossy recovery has higher rates of converging runs than FEIR and AFEIR. However, these runs are also on average much longer (up to 1000 times), and their duration has much higher variability. This is expected as with both the FEIR and AFEIR techniques several simultaneous faults may cause a failure to recover, whereas Lossy only relies on constant data and thus can always continue solving. The smaller relative differences between the expected convergence times of all methods indicate that with high page sizes, the convergence time is dominated by recoveries.

With increasing page sizes, all interpolation recoveries become less performant, however the exact ones (FEIR, AFEIR) still outperform the Lossy approach, and all these techniques remain better than the checkpointing approach for low fault rates.

8 RELATED WORK

8.1 Localizing and Detecting Errors

In a multi-node and distributed memory execution, a fail-stop error model for each process yields a localized loss of data error model globally. This naturally leads to the proposal for a Fault-Tolerant Message Passing Interface (FT-MPI), that hands control back due to a node having stopped, in order to recover data from this part of the program [17].

Chen more recently proposed to use two application specific invariant relations to detect Silent Data Corruption (SDC) [12] in CG and its derived method BiCGStab, which are the (bi-)orthogonality of search directions, and the relation between gradient and iterate. The latter consists in checking $\|b - Ax - g\| = 0$ (with our notations), which reflects the use of inherent redundancy between vectors of CG, as leveraged in this paper. Liang et al. similarly detect soft faults online in the fast Fourier transform [31].

8.2 Checkpointless Algebraic Recoveries

Many ABFT methods add checksum values such that transformations applied by a program to its data are similarly applied to these checksums. This has been accomplished for matrix-vector multiplications by adding a checksum row in a matrix [22], but also for other operations such as QR and LU factorizations [13], [19]. This approach adds little memory space overhead at the price of computational overhead. However, checks on finite precision numbers (as opposed to bitwise checks) are sensible to round-off errors, and they do not cover reduction operations.

While exploiting MPI message logging as an implicit checkpoint, Chen et al. [11] proposed an algebraic recovery method for Krylov solvers with matrix A and right-hand side b , whose iterate x is seldom passed in MPI messages.

Once all the other variables are recovered through the implicit checkpoint, the relation between the residual g and x , thus $g = b - Ax$, is used with the lost part of the iterate as unknown. This allows to recover x without checkpointing it, using inherent redundancy of the solver instead.

Langou et al. introduced the Lossy Approach with the block-Jacobi step interpolation [29]. This restart method, designed for an MPI fail-stop error model, is already extensively discussed in Section 4.3. Agullo et al. extended this work by introducing least-squares methods for the interpolation, and further studying strategies to minimize computations and communications in recoveries in the case of multiple simultaneous errors [2]. These trade-offs can naturally be applied to the interpolations we use, if needed. Agullo et al. have further extended restart recoveries relying on lost data interpolation to eigensolvers [1].

8.3 Selective Reliability

The Fault Tolerant GMRES [21] is a method consisting of GMRES iterations, run safely, enclosing a preconditioner which may run unreliably and return inexact values. Our protection of GMRES provides the complement to this selective approach tolerating errors in the preconditioner, by protecting the outer iterations – and especially the vector basis which is, according to Hoemmen et al., the most important to guarantee convergence. This work also suggests that we might not need to recover preconditioned vectors exactly, instead replacing lost data with an approximation.

9 CONCLUSIONS

This paper demonstrates that hardware DUE reporting can be exploited jointly with redundancy relations to protect iterative solvers paying a very low cost. Our two proposed methods, FEIR and AFEIR, overcome the state-of-the-art techniques in terms of overheads. They are moreover based on very simple relations that do not require deep algorithmic understanding, whereas an algorithmic technique like the Lossy Approach [29] is harder to derive. Varying page sizes up to 2MB has shown that our methods retain lower expected convergence times than the current state-of-the-art, until the edge case where vectors fit in a single memory page. At that moment, all methods become roughly equivalent due to single errors erasing full vectors. In short, our methods' efficiency is increasing with matrix size, relative to state-of-the-art methods. These straightforward low-overhead recoveries open the door to wide-spread use of algorithmic-based techniques to protect iterative methods when DUE detection is available.

Second, the paper demonstrates that by overlapping recovery with algorithmic computation, overheads can be drastically reduced. Under high error rates, of roughly more than 0.1 errors per second, the overlapping stops paying off since the chances of getting errors on non-protected computations increase, even though this trade-off is largely matrix specific. The FEIR technique provides then better performance. In any case, task-based data-flow programming models have interesting properties for resilience, not only because of inherently splitting programs into tasks, but also because overlapping computations and recoveries is done without explicit programmer intervention. Runtime support for application-level resilience could reduce the

overheads by injecting recovery tasks only when errors are encountered – this would also increase AFEIR's coverage by executing the recovery later, and still asynchronously.

Our resilience method opens the door to interesting trade-offs when SDC comes into play. Since we cover with very low overhead nearly all memory page failures, an ECC that focuses more on detecting than correcting errors would reduce SDC [27], while delegating some correction to the application level. This work will hopefully encourage future architectural, OS and runtime features to expose errors at the application level whenever lower level recoveries fail, allowing resilience aware applications to resist significantly higher fault rates than applications oblivious to resilience.

ACKNOWLEDGMENTS

This manuscript is the journal extension of a previously published conference paper [25]. This work has been partially supported by the European Research Council under the European Union's 7th FP, ERC Advanced Grant 321253, and by the Spanish Ministry of Science and Innovation under grant TIN2015-65316-P. L. Jaulmes has been partially supported by the Spanish Ministry of Education, Culture and Sports under grant FPU2013/06982. M. Moretó has been partially supported by the Spanish Ministry of Economy and Competitiveness under Juan de la Cierva postdoctoral fellowship JCI-2012-15047. M. Casas has been partially supported by the Secretary for Universities and Research of the Ministry of Economy and Knowledge of the Government of Catalonia and the Co-fund programme of the Marie Curie Actions of the European Union's 7th FP (contract 2013 BP_B 00243). We would like to thank Nicolas Vidal for his contribution on using huge pages natively.

REFERENCES

- [1] E. Agullo, L. Giraud, P. Salas, and M. Zounon, "Interpolation-Restart Strategies for Resilient Eigensolvers," *SIAM Journal on Scientific Computing*, vol. 38, no. 5, pp. C560–C583, 2016.
- [2] E. Agullo et al., "Towards resilient parallel linear Krylov solvers: recover-restart strategies," INRIA, Research Report RR-8324, 2013.
- [3] *AMD64 Architecture Programmer's Manual*. AMD, 2011, vol. 2: System Programming, ch. 9.
- [4] M. Berry et al., *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods, 2nd Edition*. Society for Industrial Mathematics, 1994.
- [5] W. Bland et al., "Extending the scope of the checkpoint-on-failure protocol for forward recovery in standard mpi," *Concurrency Computat.: Pract. Exper.*, vol. 25, no. 17, pp. 2381–2393, 2013.
- [6] M. Bougeret et al., "Checkpointing strategies for parallel jobs," in *SC*, 2011.
- [7] P. Bridges et al., "Cooperative Application/OS DRAM Fault Recovery," in *Euro-Par*, 2012, pp. 241–250.
- [8] F. Cappello et al., "Toward Exascale Resilience," *Int. J. High Perform. Comput.*, vol. 23, no. 4, pp. 374–388, 2009.
- [9] F. Cappello et al., "Towards Exascale Resilience: 2014 update," *Supercomput. front. innov.*, vol. 1, no. 1, pp. 5–28, 2014.
- [10] M. Casas, B. de Supinski, G. Bronevetsky, and M. Schulz, "Fault resilience of the algebraic multi-grid solver," in *ICS*, 2012, pp. 91–100.
- [11] Z. Chen, "Algorithm-based recovery for iterative methods without checkpointing," in *HPDC*, 2011, pp. 73–84.
- [12] Z. Chen, "Online-ABFT: an online algorithm based fault tolerance scheme for soft error detection in iterative methods," in *PPoPP*, 2013, pp. 167–176.
- [13] T. Davies and Z. Chen, "Correcting soft errors online in LU factorization," in *HPDC*, 2013, pp. 167–178.
- [14] T. Davis and Y. Hu, "The University of Florida Sparse Matrix Collection," *ACM Trans. Math. Softw.*, vol. 38, no. 1, pp. 1–25, 2011.
- [15] V. Degalahal et al., "The effect of threshold voltages on the soft error rate," in *ISQED*, 2004, pp. 503–508.

[16] A. Duran *et al.*, "OmpSs: a proposal for programming heterogeneous multi-core architectures," *Parall. Proc. Lett.*, vol. 21, no. 2, pp. 173–193, 2011.

[17] G. Fagg, A. Bukovsky, and J. Dongarra, "Fault tolerant MPI for the HARNESSE meta-computing system," in *ICCS*. Springer, 2001, pp. 355–366.

[18] K. Ferreira *et al.*, "Evaluating the viability of process replication reliability for exascale systems," in *SC*, 2011.

[19] M. Heroux *et al.*, "An Overview of the Trilinos Project," *ACM Trans. Math. Softw.*, vol. 31, no. 3, pp. 397–423, 2005.

[20] M. Heroux, J. Dongarra, and P. Luszczek, "HPCG technical specification," Sandia National Laboratory, Tech. Rep. SAND2013-8752, 2013.

[21] M. Hoemmen and M. Heroux, "Fault-tolerant iterative methods via selective reliability," Sandia National Laboratory, Tech. Rep. SAND2011-3915 C, 2011.

[22] K.-H. Huang and J. Abraham, "Algorithm-based fault tolerance for matrix operations," *IEEE Trans. Comput.*, vol. 100, no. 6, pp. 518–528, 1984.

[23] "Intel® Xeon® Processor E7 Family: Supporting next generation RAS Servers," Intel Corporation, White Paper, 2011.

[24] *Intel® 64 and IA-32 Architectures Software Developer's Manual* rev. 51. Intel Corporation, 2014, vol. 3B: System Programming Guide, Part 2, ch. 15–16.

[25] L. Jaulmes *et al.*, "Exploiting Asynchrony from Exact Forward Recovery for DUE in Iterative Solvers," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC '15. ACM, 2015, pp. 53:1–53:12.

[26] H. Kaul *et al.*, "Near-threshold voltage (NTV)-design: Opportunities and challenges," in *DAC*, 2012, pp. 1149–1154.

[27] J. Kim, M. Sullivan, and M. Erez, "Bamboo ECC: Strong, safe, and flexible codes for reliable computer memory," in *HPCA*, 2015, pp. 101–112.

[28] A. Kleen, "mcelog: memory error handling in user space," in *Linux Kongress*, 2010, pp. 159–166.

[29] J. Langou, Z. Chen, G. Bosilca, and J. Dongarra, "Recovery patterns for iterative methods in a parallel unstable environment," *SIAM J. Scientific Computing*, vol. 30, no. 1, pp. 102–116, 2007.

[30] X. Li, M. Huang, K. Shen, and L. Chu, "A realistic evaluation of memory hardware errors and software system susceptibility," in *USENIX ATC*, 2010, pp. 6–6.

[31] X. Liang *et al.*, "Correcting Soft Errors Online in Fast Fourier Transform," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC '17. ACM, 2017, pp. 30:1–30:12.

[32] A. Moody, G. Bronevetsky, K. Mohror, and B. R. d. Supinski, "Design, modeling, and evaluation of a scalable multi-level checkpointing system," in *SC*, 2010.

[33] S. Mukherjee *et al.*, "A Systematic Methodology to Compute the Architectural Vulnerability Factors for a High-Performance Microprocessor," in *MICRO*, 2003, pp. 29–42.

[34] I. Reed, "A class of multiple-error-correcting codes and the decoding scheme," *Trans. IRE Prof. Group Inf. Theory*, vol. 4, no. 4, pp. 38–49, 1954.

[35] B. Schroeder, E. Pinheiro, and W.-D. Weber, "DRAM errors in the wild: a large-scale field study," in *SIGMETRICS*, 2009, pp. 193–204.

[36] J. R. Shewchuk, "An introduction to the conjugate gradient method without the agonizing pain," Carnegie Mellon University, Tech. Rep., 1994.

[37] D. Sorin, M. Martin, M. Hill, and D. Wood, "SafetyNet: improving the availability of shared memory multiprocessors with global checkpoint/recovery," in *ISCA*, 2002, pp. 123–134.

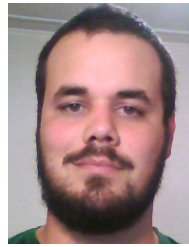
[38] O. Subasi *et al.*, "NanoCheckpoints: A Task-based Asynchronous Dataflow Framework for Efficient and Scalable Checkpoint/Restart," in *PDP*, 2015, pp. 99–102.

[39] D. Tang, P. Carruthers, Z. Totari, and M. Shapiro, "Assessment of the effect of memory page retirement on system RAS against hardware faults," in *DSN*, 2006, pp. 365–370.

[40] M. Wong *et al.*, "Towards an error model for openMP," in *Beyond Loop Level Parallelism in OpenMP: Accelerators, Tasking and More*. Springer, 2010, pp. 70–82.

[41] X. Yang, Z. Wang, J. Xue, and Y. Zhou, "The reliability wall for exascale supercomputing," *IEEE Trans. Comput.*, vol. 61, no. 6, pp. 767–779, 2012.

[42] D. H. Yoon and M. Erez, "Virtualized and Flexible ECC for Main Memory," in *ASPLOS*, 2010, pp. 397–408.



Luc Jaulmes is a PhD student at the Barcelona Supercomputing Center (BSC). He has been awarded a double degree MSc. in Engineering in computer science from Ecole Polytechnique, Paris and the Royal Institute of Technology (KTH), Stockholm, specializing in computer architecture. His focus is on the resilience of High Performance Computing (HPC) applications and memories on extreme scale machines, using novel programming models, runtime systems and runtime-aware architectures.



Marc Casas is a senior researcher at the Barcelona Supercomputing Center (BSC). He received a 5-years degree in mathematics and a PhD degree in Computer Science, from the Universitat Politècnica de Catalunya (UPC). He worked on algorithmic-based fault tolerance and active measurement methods based on software interference at the Lawrence Livermore National Laboratory (LLNL). His current research interests are high performance computing architectures, runtime systems and parallel algorithms.



Miquel Moretó received the BSc. and MSc. degrees in mathematics and electrical engineering from the Universitat Politècnica de Catalunya (UPC), and the PhD degree in the Computer Architecture Department at the same university. He spent 15 months as a postdoctoral fellow at the International Computer Science Institute (ICSI), Berkeley, USA. His research interests include studying shared resources in multithreaded architectures and modelling interconnection networks in parallel systems.



Eduard Ayguadé is a full professor in the Computer Architecture Department at the Universitat Politècnica de Catalunya (UPC), lecturing on computer organization and architecture, and on parallel programming models, since 1987. His research interests cover microarchitecture, multicore architectures, and programming models and their architectural support. He has published more than 250 papers in conferences and journals in these topics. He has participated in several research projects in the framework of the European Union and research collaborations with companies. He is associated director for research on computer sciences at the Barcelona Supercomputing Center (BSC).



Jesús Labarta is a professor at the Universitat Politècnica de Catalunya (UPC) and director of the Computer Science Department of the Barcelona Supercomputer Center (BSC). His research interests include high-performance architectures and system software, particularly programming models and performance tools. He has published more than 250 papers in conferences and journals in these topics. Labarta has a PhD in telecommunications from UPC.



Mateo Valero is a full professor in the Computer Architecture Department at the Universitat Politècnica de Catalunya (UPC), and the director of the Barcelona Supercomputing Center (BSC). His research interests focus on high performance architectures. Prof. Valero has published approximately 600 papers, served in the organization of more than 300 International conferences and given more than 400 invited talks. He has been honoured with several awards, among which the Seymour Cray, Eckert-Mauchly, Harry Goode, ACM Distinguished Service, and "King Jaime I" awards, and been named Honorary Doctor by the University of Chalmers (Sweden), Belgrade (Serbia), Las Palmas de Gran Canaria, Zaragoza and Complutense de Madrid (Spain), and Veracruz (Mexico). He is a Fellow of the IEEE, of the ACM, and an Intel Distinguished Research Fellow.

Goode, ACM Distinguished Service, and "King Jaime I" awards, and been named Honorary Doctor by the University of Chalmers (Sweden), Belgrade (Serbia), Las Palmas de Gran Canaria, Zaragoza and Complutense de Madrid (Spain), and Veracruz (Mexico). He is a Fellow of the IEEE, of the ACM, and an Intel Distinguished Research Fellow.