

1 Resiliency in Numerical Algorithm Design for 2 Extreme Scale Simulations

3 **Emmanuel Agullo**
Inria

4 **Hartwig Anzt**
KIT – Karlsruher Institut für Technologie

5 **Tommaso Benacchio**
Politecnico di Milano

6 **Hans-Joachim Bungartz**
TU München

7 **Florina M. Ciorba**
Universität Basel

8 **Daniel Drzisga**
TU München

9 **Christian Engelmann**
Oak Ridge National Laboratory

10 **Luc Giraud**
Inria

11 **Marco Heisig**
Universität Erlangen-Nürnberg

12 **Nils Kohl**
Universität Erlangen-Nürnberg

13 **Romain Lion**
University of Bordeaux

14 **Paul Mycek**
Cerfacs

15 **Enrique S. Quintana-Ortí**
Universitat Politècnica de València

16 **Ulrich Rüde**
Universität Erlangen-Nürnberg

17 **Fred Fung**
Australian National University

18 **Linda Stals**
Australian National University

19 **Samuel Thibault**
University of Bordeaux

20 **Andreas Wagner**
21 TU München

Mirco Altenbernd
Universität Stuttgart

Leonardo Bautista-Gomez
Barcelona Supercomputing Center

Luca Bonaventura
Politecnico di Milano

Sanjay Chatterjee
NVIDIA Corporation

Nathan DeBardeleben
Los Alamos National Laboratory

Sebastian Eibl
Universität Erlangen-Nürnberg

Wilfried N. Gansterer
University of Vienna

Dominik Göddeke
Universität Stuttgart

Fabienne Jézéquel
Université Paris 2 – Paris

Xiaoye Sherry Li
Lawrence Berkeley National Laboratory

Miriam Mehl
Universität Stuttgart

Michael Obersteiner
TU München

Francesco Rizzi
NexGen Analytics

Martin Schulz
TU München

Robert Speck
Jülich Supercomputing Centre

Keita Teranishi
Sandia National Laboratories – California

Dominik Thönnies
Universität Erlangen-Nürnberg

Barbara Wohlmuth
TU München



© Luc Giraud, Ulrich Rüde and Linda Stals;
licensed under Creative Commons License CC-BY

Resiliency in Numerical Algorithm Design for Extreme Scale Simulations.

Editors: John Q. Open and Joan R. Access; Article No. ; pp. :1–:57

Leibniz International Proceedings in Informatics



LIPIC Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

XX:2 Acronyms

22 Acronyms

- 23 **ABFT** *Algorithm-Based Fault Tolerance.* 22, 35
- 24 **AMR** *Adaptive Mesh Refinement.* 25, 26, 32
- 25 **API** *Application Programming Interface.* 11, 15, 16

- 26 **BLCR** *Berkeley Lab Checkpoint/Restart.* 9

- 27 **CDs** *Containment Domains.* 13, 14
- 28 **CPU** *Central Processing Unit.* 11
- 29 **CRC** *Cyclic Redundancy Checks.* 9

- 30 **DLS** *Dynamic Loop Self-scheduling.* 16
- 31 **DLS4LB** *Dynamic Loop Scheduling for Load Balancing.* 14
- 32 **DMR** *Double modular redundancy.* 20
- 33 **DRAM** *Dynamic Random-Access Memory.* 13
- 34 **DSL** *Domain Specific Languages.* 38
- 35 **DUE** *Detectable, but Uncorrectable Error.* 31

- 36 **ECC** *Error Correcting Codes.* 9, 11, 31

- 37 **FEM** *Finite Element Method.* 34
- 38 **FFT** *Fast Fourier Transforms.* 18
- 39 **FPGA** *Field-Programmable Gate Arrays.* 14

- 40 **GPU** *Graphics Processing Units.* 11, 14, 32
- 41 **GVR** *Global View Resilience.* 14

- 42 **HBM** *High-Bandwidth Memory.* 13
- 43 **HPC** *High Performance Computing.* 6, 7, 9–13, 15–17, 31, 32, 39

- 44 **LFLR** *Local-Failure Local-Recovery.* 21, 26

- 45 **MDS** *Meta Data Service.* 11
- 46 **MPI** *Message Passing Interface.* 9

- 47 **NVM** *Non-Volatile Memory.* 13

- 48 **ODE** *ordinary differential equations.* 19, 20, 35
- 49 **OS** *Operating Systems.* 9, 32

- 50 **PCG** *Preconditioned Conjugate Gradient.* 19, 22
- 51 **PDE** *Partial Differential Equations.* 19–21, 24, 34, 35
- 52 **PFS** *Parallel File System.* 8, 11, 13
- 53 **PMPI** *MPI Profiling Interface.* 11
- 54 **PVFS** *Parallel Virtual File System.* 11

- 55 **QOS** *Quality of Service.* 11

- 56 **rDLB** *robust Dynamic Load Balancing.* 15, 16

- 57 **SDC** *Silent Data Corruption.* 8, 10, 12
- 58 **SSD** *Solid-State Drives.* 13

Acronyms

XX:3

59 **TMR** *Triple Modular Redundancy.* 20

60 **ULFM** *User Level Failure Mitigation.* 13, 15

61 — **Abstract** —

62 This work is based on the seminar titled “Resiliency in Numerical Algorithm Design for Extreme
63 Scale Simulations” held March 1-6, 2020 at Schloss Dagstuhl, that was attended by all the authors.
64 Advanced supercomputing is characterized by very high computation speeds at the cost of involving
65 an enormous amount of resources and costs. A typical large-scale computation running for 48 hours
66 on a system consuming 20 MW, as predicted for exascale systems, would consume a million kWh,
67 corresponding to about 100k Euro in energy cost for executing 10^{23} floating-point operations. It is
68 clearly unacceptable to lose the whole computation if any of the several million parallel processes
69 fails during the execution. Moreover, if a single operation suffers from a bit-flip error, should the
70 whole computation be declared invalid? What about the notion of reproducibility itself: should
71 this core paradigm of science be revised and refined for results that are obtained by large scale
72 simulation?

73 Naive versions of conventional resilience techniques will not scale to the exascale regime: with a
74 main memory footprint of tens of Petabytes, synchronously writing checkpoint data all the way to
75 background storage at frequent intervals will create intolerable overheads in runtime and energy
76 consumption. Forecasts show that the mean time between failures could be lower than the time to
77 recover from such a checkpoint, so that large calculations at scale might not make any progress if
78 robust alternatives are not investigated.

79 More advanced resilience techniques must be devised. The key may lie in exploiting both
80 advanced system features as well as specific application knowledge. Research will face two essential
81 questions: (1) what are the reliability requirements for a particular computation and (2) how do we
82 best design the algorithms and software to meet these requirements? While the analysis of use cases
83 can help understand the particular reliability requirements, the construction of remedies is currently
84 wide open. One avenue would be to refine and improve on system- or application-level checkpointing
85 and rollback strategies in the case an error is detected. Developers might use fault notification
86 interfaces and flexible runtime systems to respond to node failures in an application-dependent
87 fashion. Novel numerical algorithms or more stochastic computational approaches may be required
88 to meet accuracy requirements in the face of undetectable soft errors. These ideas constituted an
89 essential topic of the seminar.

90 The goal of this Dagstuhl Seminar was to bring together a diverse group of scientists with
91 expertise in exascale computing to discuss novel ways to make applications resilient against detected
92 and undetected faults. In particular, participants explored the role that algorithms and applications
93 play in the holistic approach needed to tackle this challenge. This article gathers a broad range of
94 perspectives on the role of algorithms, applications, and systems in achieving resilience for extreme
95 scale simulations. The ultimate goal is to spark novel ideas and encourage the development of
96 concrete solutions for achieving such resilience holistically.

97 This article gathers a broad range of perspectives on the role of algorithms, applications, and
98 systems in achieving resilience for extreme scale simulations. The ultimate goal is to spark novel
99 ideas and encourage the development of concrete solutions for achieving such resilience holistically.

100 **2012 ACM Subject Classification** Theory of computation → Massively parallel algorithms; Networks
101 → Error detection and error correction; Computing methodologies → Parallel programming languages;
102 Computer systems organization → Dependable and fault-tolerant systems and networks

103 **Keywords and phrases** Numerical algorithms, Parallel computer architecture, Fault tolerance, Re-
104 siliance

105 **Digital Object Identifier** 10.4230/LIPIcs.Seminar.20101.2020.

106 **Contents**

107	1 Introduction	6
108	2 System infrastructure techniques for resilience	8
109	2.1 Detected and transparently corrected errors	9
110	2.2 Detected errors mitigated with assistance	12
111	2.2.1 Correction with incremental redesign	12
112	2.2.2 Correction with major redesign	15
113	3 Numerical algorithms for resilience	17
114	3.1 Error detecting algorithms	17
115	3.1.1 Exceptions	17
116	3.1.2 Checksums	17
117	3.1.3 Constraints	19
118	3.1.4 Technical error information	19
119	3.1.5 Multi-resolution	20
120	3.1.6 Redundancy	20
121	3.2 Error aware algorithms	21
122	3.2.1 Error aware algorithms for the solution of linear systems	21
123	3.2.2 Error aware algorithms for the solution of partial differential equations	24
124	3.3 Error oblivious algorithms	28
125	3.3.1 Gossip based methods	28
126	3.3.2 Fixed-point methods	28
127	3.3.3 Krylov subspace solvers	29
128	3.3.4 Domain decomposition	30
129	3.3.5 Time stepping	30
130	4 Future directions	30
131	4.1 Systems in support of resilient algorithms	31
132	4.1.1 Error correcting codes	31
133	4.1.2 Improving checkpoint/restart	31
134	4.1.3 Scheduler and resource management	32
135	4.2 Programming models with inherent resiliency support	33
136	4.3 Future directions for the solution of partial differential equations	33
137	4.3.1 Redundancy and replication	33
138	4.3.2 Hierarchy and mixed precision	34
139	4.3.3 Error control	35
140	4.3.4 Locality, asynchronicity and embarrassingly parallelism	35
141	4.3.5 Stochastic	36
142	4.3.6 Iterative methods	36
143	4.3.7 Low memory footprint – matrix-free	37
144	4.4 The final mile: towards a resilient ecosystem	37
145	4.4.1 Tools to support resilience software development	37
146	4.4.2 User/Programmer education	39
147	5 Conclusions	39

148 **1** Introduction

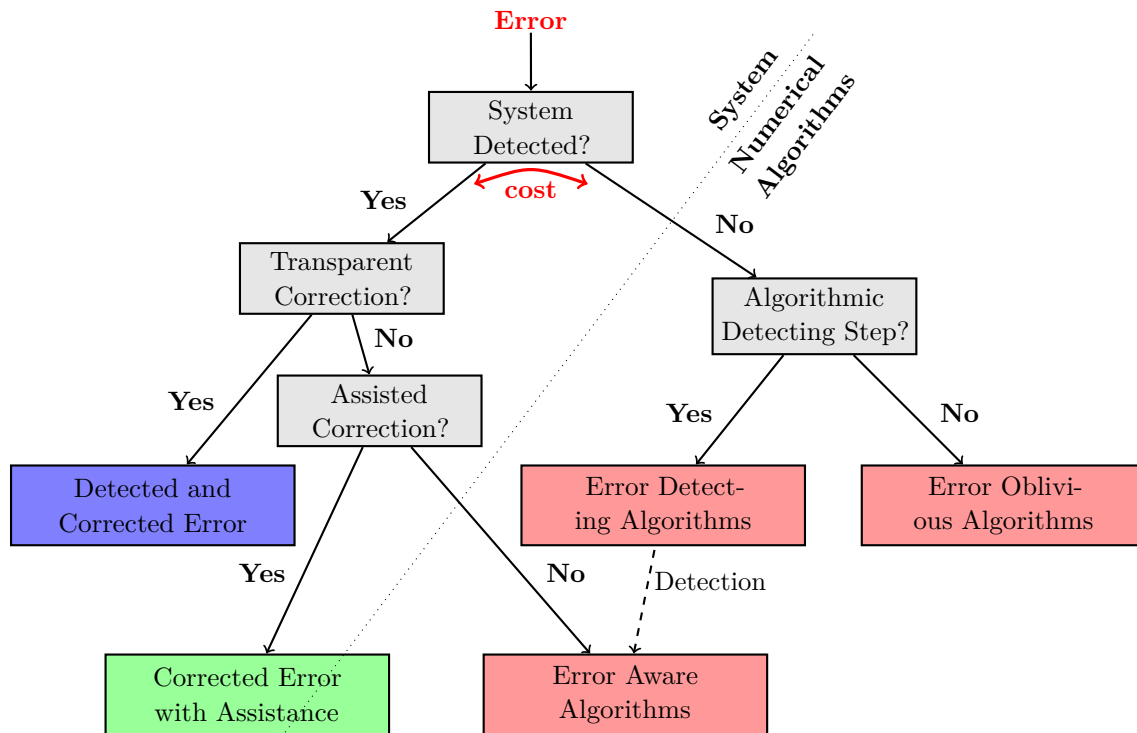
149 Numerical simulation is the third pillar in science discovery at the same level as theory
150 and experiments. To cope with the ever demanding computational resources needed by
151 complex simulations, the computational power of high performance computing systems
152 continues to increase by using an ever larger number of cores or by specialized processing.
153 On the technological side, the continuous shrinking of transistor geometry and the increasing
154 complexity of these devices affect their sensitivity to external effects and thus diminish their
155 reliability. A direct consequence is that *High Performance Computing* (HPC) applications
156 are increasingly prone to errors. Therefore the design of resilient systems and numerical
157 algorithms that are able to exploit possible unstable HPC platforms has become a major
158 concern in the computational science community. To tackle this critical challenge on the path
159 to extreme scale computation an holistic and multidisciplinary approach is required that needs
160 to involve researchers from various scientific communities ranging from the hardware/system
161 community to applied mathematics for the design of novel numerical algorithms. In this
162 article, we summarize and report on the outcomes of a Dagstuhl seminar held March 1-6,
163 2020,¹ on the topic *Resiliency in Numerical Algorithm Design for Extreme Scale Simulations*.
164 We should point out that, although error and resiliency was already quoted by J. von
165 Neumann in his first draft report on EDVAC [260, P.1, Item 1.4], it became again a central
166 concern for the HPC community in the late 2000' when the availability of the first exascale
167 computers was envisioned for the forthcoming decades. In particular, several workshops
168 were organized in the IESP (International Exascale Software Project) and EESI (European
169 Exascale Software Initiative) framework [52].

170 The hardware/system resilience community has previously defined terminology related to
171 how faults, errors, and failures occur on computing systems [18]. In this article our focus is
172 less on the cause of an error (or the underlying fault), and more on how an error presents
173 itself at the algorithmic level (or layer), impacting algorithms and applications. We thus
174 simplify the terminology often used in the hardware resilience and fault-tolerance community
175 by not using terms like soft error or hard error, and generally do not concern ourselves with
176 the reproducibility of an error (e.g., transient, intermittent or permanent). This abstraction
177 keeps the algorithmic techniques discussed herein general and applicable to a variety of fault
178 models, current architectures, and hopefully of use in future technologies.

179 To this end, we broadly categorize errors presenting themselves to the algorithmic layer
180 as either detected or undetected. Note that this categorization does not mean an error
181 is undetectable but rather that when it reached the algorithmic layer it was not detected
182 by earlier layers (e.g., hardware, operating system or middleware/system software). This
183 suggests the algorithmic layer has the opportunity to detect a previously undetected error
184 and, if possible, to deploy mitigation methods to make the algorithm resilient; effectively
185 transforming an undetected error at the algorithmic layer into a detected error. This in turn
186 may result in a failure if the algorithm is unable to handle it. For example, an undetected
187 data corruption which results in an application accessing an incorrect memory address may
188 be detectable by the algorithm but it may not be possible for the algorithm to recover what
189 the original memory address was and it may be forced to fail. If the algorithm could not
190 detect the corruption before accessing the memory region, this would conventionally end in a
191 failure (e.g., SIGSEGV issued by the operating system).

192 Many computing-intensive scientific applications that are dependent on HPC performance

¹ <https://www.dagstuhl.de/en/program/calendar/semhp/?semnr=20101>



■ **Figure 1** A classification of error handling.

193 upgrades can end up with disrupted schedules because of lack of resilience. A typical example
 194 is related to current efforts towards exascale numerical weather prediction [31, 32]. On
 195 one side, regular upgrades in weather forecast models in operations at weather centres and
 196 their spatial resolution have gone hand in hand with expanding computational resources.
 197 On the other side, scientific and socioeconomic significance of forecasts crucially hinges on
 198 tight time-bound computing schedules and timely forecast dissemination, most notably for
 199 high-impact weather events. Current disk-checkpointing schedules still take up acceptable
 200 portions of forecast runtimes, but are hardly sustainable - indeed, they already saturate
 201 file systems bandwidth. In addition, many weather forecast codes feature preconditioned
 202 iterative solvers of linear systems with several hundred thousand unknowns, many thousand
 203 times per run. Such components represent vulnerable points in a context of increasingly
 204 frequent detected and undetected errors. Novel low-overhead solutions to enhance algorithmic
 205 fault-tolerance or provide higher-level system resilience are therefore in high demand in this
 206 and other fields where nonlinear dynamics is simulated.

207 In this article we take a different approach at the classification of errors in HPC systems.
 208 In general, we try to divide errors in two main groups, those that are detected and corrected
 209 by the hardware/system (which is the focus of Section 2) and those that are detected and
 210 sometimes corrected by the numerical algorithms (Section 3). However, the HPC resilience
 211 ecosystem is not black and white, but it rather shows a wide palette of greys in between,
 212 with multiple fault tolerance tools implemented at the middleware level that are assisted by
 213 the applications/algorithms and vice-versa. Figure 1 shows this wide range of different error
 214 classifications depending on how much effort is needed at the application/algorithmic level
 215 in order to detect/correct the error.

216 The first category we observe in the leftmost leaf of the tree (blue color) is the case

217 of errors that are both detected and transparently corrected by the hardware/middleware
218 but without any intervention of the applications/algorithms. The clearest example would
219 be a detectable and correctable error in the memory generated by a single bit flip. These
220 types of errors are transparently corrected by the system without any knowledge at the
221 application/algorithmic level that such error mitigation occurred. Other examples could be
222 process replication, system-level checkpointing, process migration, among many others (see
223 Section 2.1).

224 The second category is the case of errors that are detected at the hardware/system level
225 and are mitigated at the system/middleware level (not at the algorithmic level) but with
226 assistance from the application/algorithm (green color). The most clear example of this
227 is application-based checkpointing libraries, which handle all or most of the data transfers
228 between the compute nodes and the *Parallel File System* (PFS) independently from the
229 application, but it gets hints from it to know what datasets need to be checkpointed and
230 when should the checkpoint happen. Other relevant examples are fault tolerant message
231 passing programming models and resilient asynchronous tasks. We divide these sections in
232 those approaches that require just a minor addition in the application code versus those that
233 require a complete change in the programming paradigm (see Section 2.2).

234 The other leaves of the tree (red color) correspond to those errors that cannot be corrected
235 or mitigated at the hardware/system level and have to be mitigated by changing the algorithm
236 or numerical methods to be able to tolerate those errors. We observe three different types of
237 algorithms in this branch of the tree.

238 The first type of algorithms focuses on the mitigation of errors that have been detected
239 (first red leaf from left to right), we call them error-aware algorithms. Please note that these
240 algorithms are not in charge of detecting the errors but only of mitigating them. Also, it
241 is important to notice that these algorithms do not depend on how the error was actually
242 detected; it could be hardware/middleware detection as well as algorithmic detection, in the
243 end the process of detection is irrelevant for the mitigation algorithm (see Section 3.2).

244 The second type of algorithms are those dedicated to the detection of errors that were
245 not detected at the lower levels (fourth leaf from left to right). A good example would be
246 *Silent Data Corruption* (SDC) errors that pass invisibly through the hardware but then
247 can be caught at the algorithmic level using some numerical techniques (e.g., checksum).
248 These algorithms do not try to mitigate the error per se but only detect it. Once the
249 error has been detected, it can be passed to a error aware algorithm in order to attempt a
250 correction/mitigation (see Section 3.1).

251 Finally, there also exist algorithms that can operate, tolerate and absorb errors without
252 ever being aware that the error ever occurred (last leaf to the right); we called these, error
253 oblivious algorithms. These are somehow similar to the very first (blue) category, in that the
254 errors are transparently corrected/absorbed, see Section 3.3.

255 In the following sections we discuss algorithmic and application approaches to address
256 these two categories of errors and distinguish how the approaches vary or are similar. Broadly
257 speaking, the report is divided into two parts. In Section 2 and Section 3 we discuss the
258 state-of-the-art in the areas of infrastructure and algorithms, while in Section 4 we propose
259 possible areas of interest in future research.

260 **2 System infrastructure techniques for resilience**

261 In this section we describe the state-of-the-art of hardware and system level error detection
262 and mitigation. As previously mentioned, we divide these methods in two categories, the ones

263 that mitigate the error in a completely transparent fashion, and those that require assistance
264 from the algorithmic/application level. The following subsection, Section 2.1, concentrates
265 on the methods falling in the first category. The second category is explored in Section 2.2.

266 2.1 Detected and transparently corrected errors

267 A wide range of errors can be detected and immediately corrected by various layers in the
268 system, i.e., these errors become masked or absorbed and higher level layers do not have to
269 be involved. The detection/correction mechanisms have an extra cost in terms of storage,
270 processing and energy consumption.

271 Hardware reliability

272 At the hardware level several techniques exist to detect and correct errors. Most common
273 examples are *Error Correcting Codes* (ECC) to detect and correct single bit-errors, *Cyclic*
274 *Redundancy Checks* (CRC) error correction for network packets or RAID-1 (or higher) for
275 I/O systems. A more comprehensive discussion of these features can be found in the report
276 “Towards Resilient EU HPC Systems: A Blueprint” by Radojkovic et al. [207].

277 Operating system reliability

278 *Operating Systems* (OS) have certain capabilities to interact with architectural resilience
279 features, such as ECC and machine check exceptions. OSs are mostly concerned with resource
280 management and error notification. However, some advanced OS resilience solutions exist
281 such as Mini-ckpts [92]. It is a framework that enables application survival despite the
282 occurrence of a fatal operating system failure or crash. It ensures that the critical data
283 describing a process is preserved in persistent memory prior to the failure. Following the
284 failure, the OS is rejuvenated via a warm reboot and the application continues execution
285 effectively making the failure and restart transparent. The mini-ckpts rejuvenation and
286 recovery process is measured to take 3 s to 6 s and has a failure-free overhead of 3 % to 5 %
287 for a number of key HPC workloads.

288 System-level checkpoint/restart

289 *Berkeley Lab Checkpoint/Restart* (BLCR) [125] is a system-level checkpoint/restart solution
290 that transparently saves and restores process state. In conjunction with a *Message Passing*
291 *Interface* (MPI) [95] implementation, it can transparently save and restore the process states
292 of an entire MPI application. An extension of BLCR [257, 262, 263] includes enhancements
293 in support of scalable group communication for MPI membership management, reuse of
294 network connections, transparent coordinated checkpoint scheduling, a job pause feature, and
295 full/incremental checkpointing. The transparent mechanism for job pause allows live nodes
296 to remain active and roll back to the last checkpoint, while failed nodes are dynamically
297 replaced by spares before resuming from the last checkpoint. A minimal overhead of 5.6%
298 is reported in case migration takes place, while the regular checkpoint overhead remains
299 unchanged.

300 The hybrid checkpointing technique [111] alternates between full and incremental check-
301 points: At incremental checkpoints, only data changed since the last checkpoint is captured.
302 This results in significantly reduced checkpoint sizes and overheads with only moderate
303 increases in restart overhead. After accounting for cost and savings, the benefits due to
304 incremental checkpoints are an order of magnitude larger than the overheads on restarts.

305 Silent Data Corruption (SDC) protection

306 FlipSphere [93] is a tunable, transparent *Silent Data Corruption* (SDC) detection and correc-
307 tion library for HPC applications. It offers comprehensive SDC protection for application
308 program memory using on-demand memory page integrity verification. Experimental bench-
309 marks show that it can protect 50 % to 80 % of program memory with time overheads of 7 %
310 to 55 %.

311 Proactive fault tolerance using process or virtual machine migration

312 Proactive fault tolerance [88, 187, 264] prevents compute node failures from impacting run-
313 ning applications by migrating parts of an application, i.e., tasks, processes, or virtual
314 machines, away from nodes that are about to fail. Pre-fault indicators, such as a significant
315 increase in temperature, can be used to avoid an imminent failure through anticipation and
316 reconfiguration. As computation is migrated away, application failures are avoided, which
317 is significantly more efficient than checkpoint/restart if the prediction is accurate enough.
318 The proactive fault tolerance framework consists of process and virtual machine migration,
319 scalable system monitoring and online/offline system health analysis. The process-level live
320 migration supports continued execution of applications during much of process migration
321 and is integrated into an MPI execution environment. Experiments indicate that 1 s to 6.4 s
322 of prior warning are required to successfully trigger live process migration, while similar
323 operating system virtualization mechanisms require 13 s to 24 s. This error oblivious approach
324 complements checkpoint/restart by nearly cutting the number of checkpoints by half when
325 70% of the faults are handled proactively.

326 Resiliency using task-based runtime systems

327 Task-based runtime systems have appealing intrinsic features for resiliency due to the fault
328 isolation they provide by design as they have a view of the task flow and dynamically schedule
329 task on computing units (often to minimize the time to solution or energy consumption).
330 Once an error is detected and identified by the hardware or the algorithm, the runtime system
331 can limit its propagation through the application by reasoning about the data dependencies
332 among tasks [176]. For example, one can envision the scenario where an uncorrectable
333 hardware error is detected triggering the runtime system to dynamically redistribute the
334 tasks to the remaining resources available.

335 Task-based runtime systems can also limit the size of the state needed to be saved to
336 enable restarting computations, when an error is encountered [171, 172, 254]. In classical
337 checkpoint-restart mechanisms, the size of the checkpoint can become very large for large-
338 scale applications, and managing it can take up a significant portion of the overall execution.
339 A task-based runtime system simplifies the identification of points during the application
340 execution when the state size is small, since only task boundaries need to be considered for
341 saving the state. Further, identification of idempotent tasks can greatly help task-based
342 runtimes to further reduce the overheads by completely avoiding data backups specific to
343 those tasks. Recent works on on-node task parallel programming models suggest that a
344 simple extension of the existing task-based programming framework enables efficient localized
345 recovery [200, 236, 238].

346 The checkpointing itself can also be achieved completely asynchronously [171, 172, 254].
347 The runtime allows tasks to read data being saved, and only blocks those tasks that attempt
348 to overwrite data being saved. Since the runtime system knows which data will soon be
349 overwritten by some tasks, it can prioritize the writing of the different pieces so as to have

350 as little impact on the execution as possible. At the restarting point, the runtime also has
351 all information to be able to achieve a completely local recovery. The replacement node can
352 restart from the last valid checkpoint of the previously-failed node, while the surviving nodes
353 can just replay the required data exchanges.

354 With the recent emergence of heterogeneous computing systems utilizing *Graphics Pro-*
355 *cessing Units* (GPU), the task programming model is being used to offload computation from
356 the *Central Processing Unit* (CPU) to the GPU. VOCL-FT [202] offers checkpoint/restart
357 for computation offloaded to GPU using OpenCL [118]. It transparently intercepts the com-
358 munication between the originating process and the local or remote GPU to automatically
359 recover from ECC errors experienced on the GPU during computation. Another preliminary
360 prototype design extends this concept in the context of OpenMP [42] using a novel concept for
361 *Quality of Service* (QOS) and a corresponding *Application Programming Interface* (API) [89].
362 While the programmer is specifying the resilience requirements for certain offloaded tasks,
363 the underlying programming model runtime decides on how to meet them using a QOS
364 contract, such as by employing task-based checkpoint-restart or redundancy.

365 Resilience via complete redundancy

366 The use of redundant MPI processes for error detection has been widely analyzed in the last
367 decade [55, 70, 208, 273]. Modular redundancy incurs high overhead, but offers excellent error
368 detection accuracy and coverage with few to no false positive or false negatives.

369 Complete modular redundancy is typically too expensive for actual HPC workloads.
370 However, it can make sense for certain subsystems such as parts of a PFS. The *Meta Data*
371 *Service* (MDS) of a networked PFS is a critical single point of failure. An interruption
372 of service typically results in the failure of currently running applications utilizing its file
373 system. A loss of state requires repairing the entire file system, which could take days on
374 large-scale systems, and may cause permanent loss of data. PFSs such as Lustre [67] often
375 offer some type of active/standby fail-over mechanism for the MDS. A solution [128] for the
376 MDS of the *Parallel Virtual File System* offers symmetric active/active replication using
377 virtual synchrony with an internal replication implementation. In addition to providing high
378 availability, this solution is taking advantage of the internal replication implementation by
379 load balancing MDS read requests, improving performance over the non-replicated MDS.

380 Resilience via partial redundancy

381 Partial redundancy has been studied to decrease the overhead of complete redundancy [85,
382 234, 239, 240]. Adaptive partial redundancy has also been proposed wherein a subset of
383 processes is dynamically selected for replication [108]. Partial replication (using additional
384 hardware) of selected MPI processes has been combined with prediction-based detection
385 to achieve SDC protection levels comparable with those of full duplication [37, 38, 188]. A
386 Selective Particle Replication approach for meshfree particle-based codes protects the data
387 of the entire application (as opposed to a subset) by selectively duplicating 1 % to 10 % of
388 the computations within processes incurring a 1 % to 10 % overhead [57].

389 Resilience via complete and/or partial redundancy

390 RedMPI [91] enables a transparent redundant execution of MPI applications. It sits between
391 the MPI library and the MPI application, utilizing the *MPI Profiling Interface* (PMPI) to
392 intercept MPI calls from the application and to hide all redundancy-related mechanisms. A
393 redundantly executed application runs with $r * m$ MPI processes, where m is the number of

394 MPI ranks visible to the application and r is the replication degree. RedMPI supports partial
395 replication, e.g., a degree of 2.5 instead of just 2 or 3, for tunable resilience. It also supports a
396 variety of message-based replication protocols with different consistency. Not counting in the
397 need for additional resources for redundancy, results show that the most efficient consistency
398 protocol can successfully protect HPC applications even from high SDC rates with runtime
399 overheads from 0% to 30%, compared to unprotected applications without redundancy.
400 Partial and full redundancy can also be combined with checkpoint/restart [85]. Non-linear
401 trade-offs between different levels of redundancy can be observed when additionally using
402 checkpoint/restart, since computation on non or less redundant resources is significantly less
403 reliable than computation on fully or more redundant resources.

404 **Interplay between resilience and dynamic load balancing**

405 Scheduling of application jobs at the system level contributes to exploiting parallelism by
406 placing and (dynamically) balancing the batch jobs on the local site resources. The jobs
407 within a batch are already heterogeneous; yet, current batch schedulers rarely co-allocate, and
408 most often only allocate, computing resources (while network and storage continue to be used
409 as shared resources). Dynamic system-level parallelism can arise when certain nodes become
410 unavailable (due to hard and permanent errors) or recover (following a repair operation). This
411 can be exploited during execution by increasing opportunities for system-level co-scheduling
412 in close proximity of jobs that exhibit different characteristics (e.g., co-scheduling a classical
413 compute-intensive job in close proximity to a data-intensive job) and by dynamic resource
414 reallocation to jobs that have lost resources due to failures or to waiting jobs in the queue.

415 **2.2 Detected errors mitigated with assistance**

416 In this section we focus on correction methods that need assistance from the upper layers in
417 order to achieve resilience and correctness. It is important to note that there are multiple
418 methods that offer assisted fault tolerance but some of them involve a few additional lines
419 of code while others require rewriting the whole applications using a specific programming
420 model. Therefore, we will divide this section into subsections depending on the programming
421 and/or redesign effort that is required.

422 **2.2.1 Correction with incremental redesign**

423 As explained in Section 2.1, it is possible to perform system-level checkpointing without any
424 feedback from the application or the algorithm or any upper layer. The issue with system-level
425 checkpointing is that the size (and therefore the time and energy cost) of checkpointing
426 is much larger than what is really required to perform a restart of the application. Thus,
427 application-level checkpointing is an attempt to minimize the size of checkpoints to the
428 minimum required for the application to be able to restart.

429 **Performance modeling and optimization of checkpoint-restart methods**

430 Research on simulation tools assessing the performance of certain checkpoint-restart strategies
431 is presented in various publications [15, 73, 87, 165]. Different theoretical approaches are
432 used and tools are developed that either simulate a fictional software or wrap an actual
433 application.

434 A lot of work has been done to examine and model the performance of multilevel
435 checkpointing approaches [28, 34, 156, 274]. Here, the parallel distribution of the snapshots as

436 well as the target storage system are considered as objectives for performance optimization.
437 Asynchronous techniques are considered, such as non-blocking checkpointing, where a subset
438 of processes are dedicated to manage the creation and reconstruction of snapshots [69, 219].
439 As a measure to saving storage and speeding up I/O, data compression is another subject
440 that is considered in the literature as, e.g., by Di and Cappello [74], and in one of the case
441 studies in Section 3.2.2.

442 Resilient checkpointing has been considered with the help of nonvolatile memory, as for
443 instance implemented in PapyrusKV [154], a resilient key-value blob-storage. Other resilient
444 checkpointing techniques include the self-checkpoint technique [245], which reduces common
445 redundancies while writing checkpoints, or techniques reducing the amount of required
446 memory through hierarchical checkpointing [182], or differential checkpointing [153].

447 Message logging

448 Message logging is a mechanism to log communication messages in order to allow partial restart
449 as for example examined by Cantwell et al. [51]. While improving on basic checkpointing
450 strategies, message logging-based approaches can themselves entail large overheads because
451 of log sizes. The checkpointing protocol developed by Ropars et al. [213] does not require
452 synchronization between replaying processes during recovery and limits the size of log
453 messages. Other approaches combine task-level checkpointing and message logging with
454 system-wide checkpointing [237]. This protocol features local message logging and only
455 requires the restart of failing tasks. It is also possible to combine message logging with local
456 rollback and *User Level Failure Mitigation* (ULFM) (Section 2.2.2) to improve log size [173].

457 Multilevel checkpointing libraries

458 Current HPC systems have deep storage hierarchies involving *High-Bandwidth Memory*,
459 *Dynamic Random-Access Memory*, *Non-Volatile Memory*, *Solid-State Drives* and the PFS,
460 among others. Multilevel Checkpointing libraries offer a way to leverage the different storage
461 layers in the system through a simple interface. The objective is to abstract the storage
462 hierarchy to the user, so that one does not need to manually take care of where the data
463 is stored or the multiple data movements required between storage levels. Each level of
464 checkpointing provides a different trade-off between performance and resilience, where usually
465 lower levels use close storage that offers higher performance but limited resilience, and higher
466 levels rely on stable storage (e.g., PFS), which is more resilient but slower. Mature examples
467 of multilevel checkpoint libraries are SCR [183], FTI [28], CRAFT [226] and VeloC [189].
468 Both SCR and FTI provide support via simple interfaces for storing application checkpoint
469 data on multiple levels of storage, including RAM disk, burst buffers, and the parallel file
470 system. Both SCR and FTI provide redundancy mechanisms to protect checkpoint data
471 when it is located on unreliable storage and can asynchronously transfer checkpoint data to
472 the parallel file system in the background while the application continues its execution. In
473 addition, FTI also supports transparent GPU checkpointing. Finally, VeloC is a merge of
474 the interfaces of both FTI and SRC. Note that some of these libraries offer the option for
475 keeping multiple checkpoints so that the application can roll-back to different points in the
476 past if necessary.

477 Containment Domains

478 *Containment Domains* (CDs) provide a programming construct to facilitate the preservation-
479 restoration model, including nesting control constructs, and durable storage [248]. The

480 following features are attractive for large-scale parallel applications. First, CDs respect
481 the deep machine and application hierarchies expected in exascale systems. Second, CDs
482 allow software to preserve and restore states selectively within the storage hierarchy to
483 support local recovery. This enables preservation to exploit locality of storage, rather than
484 requiring every process to recover from an error, and limits the scope of recovery to only the
485 affected processors. Third, since CDs nest, they are composable. Errors can be completely
486 encapsulated, or escalated to calling routines through a well-defined interface. We can easily
487 implement hybrid algorithms that combine both preservation-restoration and data encoding.

488 Use cases include an implementation of a parallel resilient hierarchical matrix multipli-
489 cation algorithm using a combination of ABFT (for error detection) and CDs (for error
490 recovery) [16]. It was demonstrated that the overhead for error checking and data preservation
491 using the CDs library is exceptionally small and encourages the use of frequent, fine-grained
492 error checking when using algorithm based fault tolerance.

493 **Application versioning**

494 *Global View Resilience* (GVR) [64] accommodates APIs to enable multiple versioning of
495 global arrays for the single program, multiple data programming model. The core idea is
496 the fact that naive data redundancy approaches potentially store wrong applications states
497 due to the large latency associated with error detection and notification. In addition to
498 multiple versioning, GVR provides a signaling mechanism that triggers the correction of
499 application states based on user-defined application error conditions. Use cases include an
500 implementation of resilient Krylov subspace solvers [275].

501 **Mitigating performance penalties due to resilience via dynamic load balancing**

502 Detected and corrected errors induce variation in the execution progress of applications
503 when compared to error-free executions. This can manifest itself as load imbalance. Many
504 application-level load balancing solutions have been proposed over the years and can help to
505 address this problem. We mention here a few available packages.

506 Available load balancing software includes Zoltan [252] that requires users to describe
507 the workload across processes as a graph and offers an object oriented interface. Further we
508 mention *Dynamic Loop Scheduling for Load Balancing* (DLS4LB) [54], a recently developed
509 library for MPI applications that contains a portfolio of self-scheduling based algorithms for
510 load balancing. StarPU [254] proposes support for asynchronous load-balancing [172] for
511 task-based applications. The principle is to let the application submit only a part of its task
512 graph, let some of it execute on the platform and observe the resulting computation balance.
513 A new workload distribution can then be computed and the application is allowed to submit
514 more of the task graph, whose execution can be observed as well. OmpSs [80] is an effort to
515 extend OpenMP in order to support asynchronous execution of tasks including a transparent
516 interface for hardware accelerators such as GPUs and FPGAs. OmpSs is built on top of the
517 Mercurium compiler [251] and the nanos++ runtime system [249].

518 HCLib [271] is a task-based programming model that implements locality-aware runtime
519 and work-stealing. It offers a C and C++ interface and can be coupled with inter-process
520 communication models, such as MPI. Charm++ [151] features an automatic hierarchical
521 dynamic load balancing method that overcomes the scalability limitation of centralized
522 load balancing as well as the poor performance of completely distributed systems. Such
523 a technique can be triggered dynamically after a failure hits the system and the workload
524 needs to be redistributed across workers.

525 2.2.2 Correction with major redesign

526 The correction of some detected errors might have a strong impact of the algorithm that
527 has to implement the mitigation. The mitigation design can be made more affordable if
528 some components of the software stack have already some appealing features to handle such
529 situations.

530 Resilience support in the Message Passing Interface (MPI)

531 Most MPI implementations by default are designed to terminate all processes when errors are
532 detected. However, this termination occurs irrespective of the scope of the error, requiring
533 global shut-down and restart even for local errors in a single process. This inherent scalability
534 issue can be mitigated if MPI keeps all survived processes to continue and/or if restart
535 overheads are reduced. The MPI community has proposed several recovery approaches, such
536 as FA-MPI [127] or MPI-ULFM [41] to enable alternatives of global shut-down, as well as
537 better error handling extensions, like MPI_Reinit [159], to reduce overhead and impact
538 of failures. Among these approaches, MPI-ULFM is the most advanced and well known.
539 It provides a flexible low-level API that allows application specific recovery via new error
540 handling approaches and dynamic MPI communicator modification under process failures,
541 although with significant complexities for the application developer using the new APIs.
542 Several approaches have been proposed to mitigate this complexity by creating another set
543 of library APIs built atop of MPI-ULFM [51, 99, 100, 225, 253]. However, as of now, in part
544 due to its complexity when used on real-world applications and limited support in system
545 software, MPI-ULFM as a whole has not been adopted in the MPI standard and hence is
546 not readily usable for typical HPC application programmers. Nevertheless, various aspects
547 of ULFM are in the process of standardization and will provide more mechanisms in MPI
548 to build at least certain fault tolerant applications, starting with the upcoming MPI 4.0
549 standard.

550 Resilience abstractions for data-parallel loops

551 Data-parallel loops are widely encountered in N -body simulations, computational fluid
552 dynamics, particle hydrodynamics, etc. Optimizing the execution and performance of such
553 loops has been the focus of a large body of work involving dynamic scheduling and load
554 balancing. Maintaining the performance of applications with data-parallel loops running in
555 computing environments prone to errors and failures is a major challenge. Most self-scheduling
556 approaches do not consider fault-tolerance or depend on error and failure detection and react
557 by rescheduling failed loop iterations (also referred to as tasks). A study of resilience in
558 self-scheduling of data-parallel loops has been performed using SimGrid-based simulations of
559 highly unpredictable execution conditions involving various problem sizes, system sizes, and
560 application and systemic characteristics (namely, permanent node failures), that result in
561 load imbalance [241]. Upon detecting a failed node, re-execution is employed to reschedule
562 the loop iterations assigned to the failed node.

563 A *robust Dynamic Load Balancing* (rDLB) approach has recently been proposed for
564 the robust self-scheduling of independent tasks [180]. The rDLB approach proactively and
565 selectively duplicates the execution of assigned chunks of loop iterations and does not depend
566 on failure or perturbation detection. For exponentially distributed permanent node failures,
567 a theoretical analysis shows that rDLB is linearly scalable and its cost decreases quadrati-
568 cally with increasing system size. The reason is that increasing the number of processors
569 increases the opportunities for selectively and proactively duplicating loop iterations to

570 achieve resilience. rDLB is integrated into a dynamic loop scheduling library (DLS4LB, see
571 Section 2.2.1) for MPI applications. rDLB enables the tolerance of up to $(P - 1)$ process
572 failures, where P is the number of processes executing an application. For execution environ-
573 ments with performance-related fluctuations, rDLB boosts the robustness of *Dynamic Loop*
574 *Self-scheduling* (DLS) techniques by a factor up to 30 and decreases application execution
575 time up to 7 times compared to their counterparts without rDLB.

576 Resilience extension for performance portable programming abstractions

577 With the increasing diversity of the node architecture of HPC systems, performance portability
578 has become an important property to support a variety of computing platforms with the
579 same source code while achieving a comparative performance to those programmed with the
580 platform specific programming models. Today, Kokkos [82] and Raja [30,250] accommodate
581 modern C++ APIs to permit an abstraction of data allocation and parallel loop execution for
582 a variety of runtime software and node architectures. This idea can be extended to express
583 the redundancy of data and computation to achieve resilience while hiding the details of
584 the data persistence and redundant computation. Recently, the resilient version of Kokkos
585 was proposed for a natural API extension of Kokkos' data (memory space) and parallel
586 loop (execution space) abstractions to (1) enable resilience with minimal code refactoring
587 for the applications already written with Kokkos and (2) provide common interface to call
588 any external resilience libraries such as VeloC [189]. The new software will be released in a
589 special branch in <https://github.com/kokkos/kokkos>.

590 The resilience abstraction idea has also been applied to task parallel programming models
591 such as Charm++ [151], HCLib [271], HPX [150], OmpSs [80] and StarPU [254] to integrate a
592 variety of resilient task program execution options such as replay, replication, algorithm-based
593 fault tolerance and task-based checkpointing. Task-based programming models indeed have
594 a very rich view over the structure of the application computation, and notably its data,
595 and have a lot of control over the computation execution, without any need for intervention
596 from the application. Replaying a failed task consists of issuing it again with the same input,
597 discarding the previous erroneous output, and replicating a task consists of issuing it several
598 times with different output buffers and comparing the result. Dynamic runtime systems can
599 then seamlessly introduce replay and replication heuristics, such as trying to run different
600 implementations and/or computation units, without the application having to be involved
601 beyond optionally providing different implementations to be tried for the same task.

602 The task graph view also allows for very optimized checkpointing [171,172,254]. In the
603 task-based programming model, each checkpoint is a cut in the task graph, which can be
604 expressed trivially within the task submission code, and only the data of the crossing edges
605 need to be saved. Even better, the synchronization between the management of checkpoint
606 data and application execution can be greatly relaxed. The transfer of the data to the
607 checkpoint storage can indeed be started as soon as the data is produced within the task
608 graph, and not only once all tasks before the checkpoint are complete. A checkpoint is
609 then considered complete when all its pieces of data have been collected. It is possible that
610 tasks occurring after the checkpoint may run to completion before the checkpoint itself is
611 completed. All in all, this allows for a lot more time for the data transfers to complete, and
612 lessens the I/O bandwidth pressure.

613 **Software engineering approaches for resilience by design**

614 Resilience design patterns [141,142] offer an approach for improving resilience in extreme-scale
615 HPC systems. Frequently used in computer engineering, design patterns identify problems and
616 provide generalized solutions through reusable templates. Reusable programming templates
617 of these patterns can offer resilience portability across different HPC system architectures
618 and permit design space exploration and adaptation to different (performance, resilience,
619 and power consumption) design trade-offs. An early prototype [14] offers multi-resilience for
620 detection, containment and mitigation of silent data corruption and MPI process failures.

621 **3 Numerical algorithms for resilience**

622 In this section, we focus on the handling of errors at the algorithmic level. We see three
623 different classes of problems to tackle here: (i) detection of un-signaled errors (mostly bit
624 flips and other instances of silent data corruption, Section 3.1), (ii) correction of errors that
625 have been signaled but could not be corrected at the hardware or middleware layer (by error
626 aware algorithms, Section 3.2), (iii) design of error oblivious algorithms that deliver the
627 correct result even in the presence of (not too frequent) errors (Section 3.3).

628 In addition to correctness in the presence of errors, an important challenge in all our
629 considerations is efficiency in terms of algorithm runtime. In this context, additional
630 algorithmic components such as work stealing and asynchronous methods (where missing
631 data are simply an extreme case of delay) have to be considered. We mention these methods
632 when describing methods that can make use of such runtime optimizing measures.

633 **3.1 Error detecting algorithms**

634 In this section, we focus on mechanisms to numerically detect errors that have not been
635 detected by the underlying system or middleware. We have identified several techniques that
636 allow us to (likely) notice the occurrence of an error at several layers of numerical algorithms.
637 Table 1 gives an overview of some detection techniques and the algorithmic components or
638 numerical methods where they are applicable.

639 **3.1.1 Exceptions**

640 Exceptions are a way a program signals that something went wrong during execution. We
641 consider the case where exceptions are caused by data corruption that can, for example, lead
642 to division by zero or out-of-range access. Most programming languages support a way of
643 handling exceptions. The algorithm programmer can register an exception handler that gets
644 called whenever an exception occurs. If the error is recoverable, the exception handler will
645 specify how best to continue afterwards. If the error is not recoverable, the program will be
646 aborted. Exceptions are a straight-forward way to detect certain types of errors and can be
647 applied to all numerical algorithms. However, they obviously only see a small subset of all
648 possible errors and it is not trivial to decide when to use exceptions handlers in the light of a
649 trade-off between correctness, robustness and runtime efficiency.

650 **3.1.2 Checksums**

651 Checksums could be used at the hardware or middleware layer to detect errors, but here we
652 will discuss checksums as employed on the algorithmic layer where we have a more detailed

■ **Table 1** Numerical error detection: Overview of error detection techniques and numerical ingredients and methods where they are applied. Note that we mark a method as applicable only if it is or can be used in the respective algorithm itself, not only at lower level functionality, i.e., we do not mark checksums for multigrid as checksums are only used in the BLAS 2/3 kernels used as inner loops or in the GS/J/SOR smoothers.

	exceptions	checksum	constraints	tech error	multi resolution	redundancy
BLAS 2/3	×	×				×
Direct Solvers	×	×				×
Krylov	×		×	×		×
Multilevel / Multigrid	×			×	×	×
Domain Decomposition	×					×
GS/Jac/SOR	×			×		×
Nonlinear Systems	×			×		×
Time Stepping (ODEs)	×			×	(×)	×
PDEs	×	×	×	×	×	×
Quadrature	×		×	×	×	×

653 knowledge about the existence of numerical or algorithmic invariants. Checksum techniques
 654 have been used in various numerical algorithms. We list some examples below.

655 **BLAS 2/3:** Checksum encoding matrices, introduced by Huang and Abraham [137] requires
 656 (i) adding redundant data in some form (encoding), (ii) redesign of the algorithm to operate
 657 on the respective data structures (processing), and (iii) checking the encoded data for errors
 658 (detection). We ignore the recovery phase here and refer to Section 3.2. Checksums are used
 659 in FT-ScaLAPACK [267] for dense matrix operations such as MM, LU and QR factorization
 660 and more recently in hierarchical matrix multiplication [16]. Wu et al. give a good survey of
 661 checksum deployment in dense linear algebra [268].

662 **Gauss-Seidel/Jacobi/SOR and multigrid:** In [179], checksums are used to detect errors
 663 in the Jacobi smoother, the restriction and interpolation operators of a multigrid method
 664 solving a two-dimensional Poisson equation.

665 **Krylov subspace methods:** Tao et al. propose a new checksum scheme using multiple
 666 checksum vectors for sparse matrix-vector multiplication, which is shown to be generally
 667 effective for several preconditioned Krylov iterative algorithms [247]. Also [1, 227] use
 668 checksums for protection within the conjugate gradient (CG) algorithm.

669 **FFT:** Checksum can also be used in *Fast Fourier Transforms* (FFT)s similarly as in
 670 matrix-vector multiplication. Liang et al. [167] develop a new hierarchical checksum scheme
 671 by exploiting the special discrete Fourier transform matrix structure and employ special
 672 checksum vectors. Checksums are applicable to many important kernels such as matrix-
 673 matrix multiplication, but are costly. In addition, it can be difficult to specify a suitable
 674 threshold for ‘equality’ in the presence of round-off errors. For many numerical calculations
 675 such as scalar products, checksums are not applicable at all.

676 3.1.3 Constraints

677 In some applications, constraints for different types of variables are known. Examples are
678 positivity constraints, conservation laws for physical quantities or known bounds for internal
679 numerical variables.

680 **Krylov subspace methods:** Resilience was already of importance in the early days of
681 digital computers. In the original PCG paper [132], Hestenes and Stiefel noticed that the
682 reciprocal value of α (the step length) is bounded above (repectively, below) by the reciprocal
683 of the smallest eigenvalues (respectively the inverse of the largest eigenvalue) of the matrix.
684 The inequality involving the largest eigenvalue (for which in practice it may be cheaper to
685 get an approximation) was used to equip PCG with error detection capabilities in [1].

686 **Partial differential equations:** Checking for bounds can be associated with minimal or
687 extremely high cost depending on whether extra information has to be computed (such as
688 eigenvalues of matrices) or not. Reliability is, in general, an issue as only those errors leading
689 to violation of these constraints can be detected. An example of the use of problem-informed
690 constraints can be found in [186]. In this work, the authors derive a priori bounds for the
691 discrete solution of second-order elliptic PDEs in a domain decomposition setting. Specifically,
692 they show that the bounds take into account the boundary conditions, are cheap to compute,
693 general enough to apply to a wide variety of numerical methods such as finite elements
694 or finite differences, and provide an effective way to handle faulty solutions synthetically
695 generated.

696 3.1.4 Technical error information

697 In many numerical large scale applications, the main computational task involves the
698 approximate computation of integrals, algebraic systems, systems of ODEs or PDEs. For
699 all these problems, various types of error information such as residuals, differences between
700 iterations, round-off error estimates and discretization error estimates can be used as indicators
701 of errors either by their size or by monotonicity criteria. We give several examples from
702 literature for different classes of numerical algorithms.

703 **Krylov subspace methods:** Round-off error bounds can be used in Krylov subspace
704 methods. They fit in the general framework of round-off error analysis [133] and have been
705 considered in the context of Krylov subspace methods in finite precision arithmetic [169, 178].

706 Vorst and Ye proposed a residual gap bound [256] (bound for the norm of the residual
707 gap between the true and the computed residuals) based on round-off error analysis that was
708 later used as a criterion for actual error detection in [1] when bit flips occur. The detection
709 of errors in Krylov methods via violation of orthogonality is proposed in [63].

710 **Multigrid:** Calhoun et. al [50] apply a residual/energy norm-based error detection for
711 algebraic multigrid. They use two criteria: (i) the reduction of the residual norm as a weak
712 criterion and (ii) the reduction of the quadratic form

$$713 \quad E(\mathbf{x}) = \langle A\mathbf{x}, \mathbf{x} \rangle - 2\langle \mathbf{x}, \mathbf{b} \rangle,$$

714 when solving the linear system $A\mathbf{x} = \mathbf{b}$ for symmetric positive matrices.

715 The quadratic for E calculated at level i during the down-pass of a V-cycle should be
716 less than the energy calculated at level i during the down-pass of the next V-cycle.

717 When using the full approximation scheme residual norm reductions can also be verified
718 at each level in the hierarchy of a multigrid-cycle. The structure of the full approxima-
719 tion scheme additionally provides smart recovery techniques utilizing its lower resolution
720 approximations [11].

721 **Time-stepping:** For iterative time-stepping with spectral deferred corrections, monitoring
 722 the residual of the iteration can be used to detect errors in the solution vectors [119]. In
 723 the context of parallel-in-time integration with parareal, consecutive iterates are considered
 724 in [191] to detect errors in the solution vector. In [35], an auxiliary checking scheme in
 725 contrast to the original base scheme is used to detect and correct errors during implicit and
 726 explicit time-integration. Estimating the local truncation error with two different methods
 727 is used in [121] to implement a resilient, high-order Runge-Kutta method. This “Hot Rod”
 728 approach is then also used for error correction.

729 3.1.5 Multi-resolution

730 Multi-resolution means that information is available at different resolution levels, in terms of
 731 spatial discretization (PDE), time discretization (ODE and PDE), order of discretization
 732 (PDE in space and time), matrix dimensions (numerical linear algebra, multigrid), frequencies,
 733 and so on. This leads to a certain redundancy – not an artificially introduced, but an inherently
 734 available one. This redundancy can be used to detect discrepancies or anomalies and, hence,
 735 errors that could not be detected by the system. There are numerous examples for the
 736 mentioned problem classes, we outline one example in more detail here.

737 **Sparse grids / Combination technique:** Sparse grids [48] are one particular class
 738 of multi-resolution methods. There, via the use of hierarchical bases, certain structures
 739 often seen in d -dimensional data can be exploited to alleviate the curse of dimensionality,
 740 without a significant loss of accuracy. Sparse grids have been successfully used in a wide
 741 range of problem classes where spatial discretization plays a role, such as interpolation [145],
 742 quadrature [47, 109, 110], solvers for PDEs [124, 129], or machine learning tasks [104, 105, 201]
 743 (e.g., classification, regression, clustering, or density estimation). One particular incarnation
 744 of sparse grid methods is the so-called combination technique [117]. There, based on an
 745 extrapolation-style approach, a linear combination of a specific set of full, but very coarse-grid
 746 solutions is used to get a sparse fine-grid solution. The various coarse grid solutions can be
 747 obtained in a completely independent way, using (parallel) standard solvers. This opens the
 748 way to (1) a natural two-level parallelization and to (2) an easy and cheap detection of system
 749 undetected errors: Since we actually compute solutions for the same problem on different
 750 (i.e., differently discretized) grids anyway, we can use these to detect anomalies – just by
 751 comparing the available solutions. And the detection leads immediately to a mitigation
 752 strategy (see Section 3.2.2), since we can easily exchange coarse grids in case of errors, just
 753 by changing the combination pattern [8, 9, 123, 130, 134, 192]. Therefore, this is an example
 754 for a smart algorithm that is able to do both detection and mitigation.

755 Further examples are mentioned in Section 3.1.4 as multi-resolution typically comes with
 756 corresponding error estimates based on differences between solutions at different resolution
 757 levels: multigrid and parallel time stepping.

758 3.1.6 Redundancy

759 Redundancy is a strategy for error detection that can be applied to all of the numerical algo-
 760 rithms mentioned in Table 1. It covers two approaches. In the first approach computational
 761 resources may be replicated twice or thrice. Such instances are called DMR [144, 265] or
 762 TMR [223, 261]. In the second approach the computations are repeated twice or thrice on the
 763 same resource [17, 259]. An advantage of this approach is the flexibility at the application
 764 level. Note that the first approach costs more in space or resources, the second approach
 765 costs more in time.

766 The redundancy based error detection technique described in [33] relies on in-depth
767 analysis of application and platform dependent parameters (such as the number of processors
768 and checkpointing time) to formalise the process of both resource and computation replication.
769 It provides a closed-form formula for optimal period size, resource usage and overall efficiency.

770 Ainsworth et. al [5] use replication of fault-prone components as an error detection
771 technique in a multigrid method. Also error detection in the time stepping methods from [35]
772 mentioned in Section 3.1.4 can be interpreted as redundancy based error detection.

773 The main disadvantage of replication is its cost in terms of performance, although recom-
774 puting only some instructions instead of the whole application lowers the time redundancy
775 overhead [193]. However, redundancy in some calculations should in particular be considered
776 as a possible strategy for error detection as in modern supercomputers the cost of arithmetic
777 operations tends to decrease compared to communication time.

778 3.2 Error aware algorithms

779 In this section, we look at error correction techniques within an application. We assume that
780 the application has been notified that part of the algorithm's data is corrupted or lost. In
781 that context, mitigation or containment actions have to be undertaken at the algorithmic
782 design level, where the appropriate actions depend on the data detection granularity and
783 how the notification mechanism was activated. It is possible to design both lossy and lossless
784 mitigation procedures that are tailored to the numerical algorithms under consideration.

785 In Section 3.2.1 we give a brief literature overview of ideas that can be used to complement
786 numerical mitigation or containment procedures. Then, in Section 3.2.2 we offer a more
787 detailed discussion of some recent successful attempts by presenting a few case studies in the
788 context of the solution of *Partial Differential Equations* (PDE).

789 3.2.1 Error aware algorithms for the solution of linear systems

790 A wealth of literature already exists on various, mostly isolated ideas and approaches that have
791 appeared over time. Checkpoint-restart methods are the most generic approaches towards
792 resilience for a broad spectrum of applications, see Section 2.2.1 for an introduction. We
793 first describe a general mental model to design resilient numerical algorithms independent of
794 actual machine specifications that lead to what is nowadays referred to as *Local-Failure Local-*
795 *Recovery* (LFLR) techniques. Then we move to 'classical' algorithm-based fault tolerance,
796 which originally was developed to detect and correct single bit flips on systolic architectures
797 devoted to basic matrix computations, see Section 3.1.2. Finally, we discuss a range of ideas
798 and techniques not covered by the case studies below.

799 Local-failure local-recovery

800 As far back as a decade ago, an abstract framework was developed to separate algorithm
801 design from unclear machine specifications, see also Section 2.2.1. The idea of a selective
802 reliability model as introduced by Hoemmen [45,135] is machine-oblivious and highly suitable
803 for algorithm design for machines with different levels of (memory) reliability. It has led to
804 the concept of *Local-Failure Local-Recovery* (LFLR) [253]. This model provides application
805 developers with the ability to recover locally and continue application execution when a
806 process is lost. In [253], Teranishi and Heroux have implemented this framework on top of
807 MPI-ULFM (Section 2.2.2) and analyzed its performance when a failure occurs during the
808 solution of a linear system of equations.

809 **Original algorithm-based fault tolerance with checksums**

810 The term *Algorithm-Based Fault Tolerance* (ABFT) was originally coined in conjunction
811 with protecting matrix operations with checksums to handle bit flips [136], mostly assuming
812 exact arithmetic calculation for detection and mitigation. (See Section 3.1.2 for a more
813 detailed discussion on checksums). The main drawback of checksums is that only limited
814 error patterns can be corrected and its robust practical implementation in finite precision
815 arithmetic can be complicated to tune to account for round-off errors. A second drawback
816 is that the checksum encoding, detection and recovery methods are specific to a particular
817 calculation. A new scheme needs to be designed and proved mathematically for each new
818 operation. A further drawback is to tolerate more errors, more encoded data is needed, which
819 may be costly both in memory and in computing time.

820 ABFT concepts have been extended to process failures for a wide range of matrix
821 operations both for detection and mitigation purposes [44, 62, 79, 147, 269] and general
822 communication patterns [149]. ABFT has also recently been proposed for parallel stencil-
823 based operations to accurately detect and correct silent data corruptions [58]. In these
824 scenarios the general strategy is a combination of checkpointing and replication of checksums.
825 In-memory checkpointing [147] can be used to improve the performance. The main advantage
826 of these methods is their low overhead and high scalability.

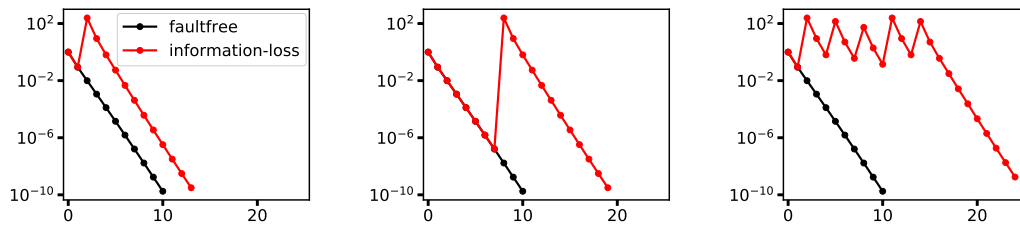
827 In practice, the significance of a bit flip strongly depends on its location, i.e., which bit in
828 the floating point representation is affected. Classical ABFT has been extended to take into
829 account floating point effects in the fault detection (checksums in finite precision) as well
830 as in the fault correction and to recover from undetected errors (bit flips) in all positions
831 without additional overhead [181].

832 **Iterative linear solvers**

833 Iterative linear solvers based on fixed point iteration schemes are, in general, examples of
834 error oblivious algorithms, as described in Section 3.3. The convergence history of the scaled
835 residual norm observed within the iterative scheme often resembles the curves displayed in
836 Figure 2. In this case the iterative scheme is a multigrid method, as in [115, 138]. The peaks
837 in the residual occur after data has been lost and when the iterations are allowed to restart
838 with some form of replacement of the lost data. In the simplest case, the lost data may just
839 be re-initialized with the value of zero, and recovery techniques to obtain better solutions
840 are discussed in Section 3.2.2.

841 It can be seen that, depending on when in the course of the iteration a small portion of
842 the approximate solution suffers from an error, we observe a delay in convergence, directly
843 proportional to an increase in runtime. In the case where errors appear too often, the solver
844 might not recover and other mitigation actions might have to be considered.

845 Explicit recovery at the algorithmic level from undetected errors have been studied
846 for iterative linear solvers [195]. In contrast to restarting, a number of algorithm based
847 recovery strategies have been proposed, including approximate or heuristic interpolation
848 methods [2]. An approach of exactly recovering the state of the iterative solver before the
849 node failure has been investigated for the *Preconditioned Conjugate Gradient* (PCG) and
850 related methods [164, 196]. This also includes studying scenarios with multiple simultaneous
851 node failures [197] and scenarios where no replacement nodes are available [194].



■ **Figure 2** Convergence history of the residual norm as a function of the iteration count for three examples of information loss. From left to right: early, late, and multiple times

852 Approximated recovery and restart in sparse numerical linear algebra

853 For matrix computations, eigensolvers or basic kernels such as iterative linear system solvers,
 854 some recovery ideas rely on forming a small dimensional linear algebra problem where
 855 the inputs are the still valid data and the unknowns are the lost/corrupted ones. The
 856 outcome of this procedure is subsequently used to replace the lost/corrupted data and the
 857 numerical algorithm is somehow started again from that meaningful initial guess. The
 858 recovery procedure is tailored to the actual numerical algorithm. As an example, consider a
 859 fixed point iteration scheme for a linear system and suppose the lost data are entries of the
 860 iterate vector, the most dynamically evolving data in this computational scheme. Matrix
 861 entries of the iteration scheme related to the lost data, as well as some neighbouring entries,
 862 serve to build the left-hand side of a linear problem (either a linear system or a least-square
 863 problem) while the right-hand side is built from valid data. The solution of this small problem
 864 is then used to replace the corresponding lost entries of the iterate vector. The complete,
 865 updated vector is taken as a new initial guess when restarting the fixed point iteration. If
 866 the data is not corrupted too often the classical convergence theory still applies and because
 867 the new initial guess incorporates updates from the calculations performed before the error
 868 was detected, the global convergence rate is not strongly affected. The method described
 869 in adaptive recovery techniques for extreme scale multigrid in Section 3.2.2 is an example
 870 application of this technique.

871 For numerical schemes based on nested subspace search, such as Krylov subspace methods,
 872 closely related techniques have been successfully applied both for eigensolvers and linear
 873 solvers that further exploit the sparsity structure of the matrices to reduce the computational
 874 cost associated with the recovery procedure. At the cost of a light checkpoint performed
 875 once when starting the linear solver (mostly the matrix and the right-hand side vector in
 876 case of linear system solution) this mitigation approach has no overhead if the data is not
 877 corrupted during the solution computation. We refer to [2, 3, 161] for some illustrations on
 878 those numerical remedies in a parallel distributed memory framework and to [146] where
 879 these ideas are exploited for a lower granularity of data loss in a task-based runtime system.
 880 See Section 2 for references relevant to task-based runtime systems.

881 We also note that these ideas can be extended to hybrid iterative/direct numerical
 882 schemes, that have a domain decomposition flavor, where the recovery procedure can be
 883 enriched with additional features of the parallel numerical scheme such as redundancy or
 884 properties of the preconditioners [4]. They can also be extended to the time domain in the
 885 context of multilevel parallel-in-time integration techniques [229].

886 3.2.2 Error aware algorithms for the solution of partial differential 887 equations

888 The ideas introduced above in Section 3.2.1 are application agnostic but naturally apply to
889 linear systems arising from the discretization of a PDE. In that latter case, more information
890 from the underlying PDE can be closely tailored to intrinsic features of solvers such as
891 multigrid. In this section we discuss some research works on mitigation and containment
892 that exploit the properties of PDEs to aid the recovery techniques. We also present some
893 mitigation processes that are only relevant in the PDE setting.

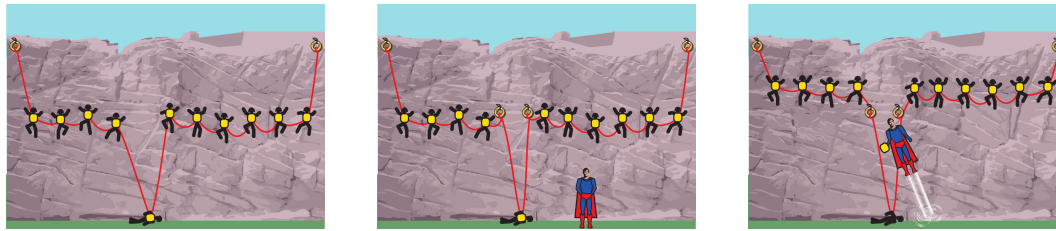
894 Adaptive recovery techniques for extreme scale multigrid

895 Some of the most efficient solvers of PDE, such as parallel geometric multigrid methods
896 [114, 143], can be based on the exchange of ghost layers in a non-overlapping domain
897 partitioning. This automatically leads to a redundancy in interface data between subdomains
898 that in turn permits the design of an efficient two-step recovery strategy for iterative solvers.
899 This is of particular interest in large-scale parallel computations. When each subdomain
900 is large, then the ratio between the data on its surface and the volume data in its interior
901 becomes small.

902 When a processor fails, the information within one or several subdomains is lost. For
903 the recovery and continued solution, the redundant ghost layer information is used in a
904 first step, to recover locally either Dirichlet- or Neumann-type data for the subdomains.
905 The global problem can then be formulated in two partitions, the outer healthy subdomain
906 and the inner faulty subdomain, where the recovery must reconstruct the lost data. Both
907 subproblems must be bi-directionally coupled via the interface and the corresponding ghost
908 layers of unknowns.

909 After re-initialization, the corrupted and reinitialized data could pollute the solution
910 globally, meaning that the locally increased error in the faulty domain can spread globally and
911 thus also affect the healthy subdomain. In order to avoid this pollution, the communication
912 between the healthy and faulty sub-problems is interrupted during the second step of the
913 recovery process. In the second step, we continue with the original iterative solver restricted
914 to the healthy sub-problem and select a suitable one for the faulty one. After some number
915 of asynchronous iteration steps both sub-problems are reconnected, see [138], and the global
916 iterative solver strategy is resumed. Note that the reconnecting step is mandatory for the
917 convergence of the iterative solver. The tearing step separating the subdomains is mandatory
918 to preserve the accuracy of the dynamic data in the healthy sub-problem, and without this
919 step the corrupted data from the faulty sub-domain pollutes the global solution. Of critical
920 importance for the performance of the method are the accuracy of the faulty sub-problem
921 solver at re-connection time and the time spent in the recovery mode. In the faulty domain,
922 the lost data can be initialized with 0, or, alternatively, compressed checkpointed data can be
923 used as described in the following section on compression techniques for checkpoint-restart.
924 Note, however, that with straight-forward compression techniques, compressed checkpoint
925 data will only be useful to recover the low frequency components in the faulty domain. If the
926 local recovery is performed with multigrid, then the low frequencies are in any case cheap to
927 recover, as compared to the cost of recomputing the lost high frequency components.

928 The accuracy within a multigrid strategy can be easily controlled by a hierarchical sum
929 of weighted residuals [216]. The overhead cost for the a-posteriori error indicator is quite
930 small compared to the overall solver cost. Having an estimate for the algebraic error in both
931 sub-problems at hand, the re-connection step is determined automatically. To speed up the



■ **Figure 3** Illustration of the steps in the adaptive recovery technique for extreme scale multigrid. Left: A detectable error occurred. Middle: The communication between the healthy and faulty sub-domains is interrupted. Right: The original iterative solver restricted to the healthy domain continues while another suitable solver is asynchronously used in the faulty domain. Once the solution in the faulty domain reaches a certain accuracy, the communication between the domains is re-enabled.

932 time which is spent in the recovery, a so-called ‘superman strategy’ is applied [138], see also
 933 Figure 3 for an illustration. More resources compared to the situation before the fault are
 934 allocated to the faulty sub-problem. A short recovery phase in combination with carefully
 935 selected re-coupling criteria then guarantees a highly efficient fault-tolerant solver.

936 Of special interest is a massively parallel multigrid method as base solver. In combination
 937 with the tearing and intersection approach for the recovery, it results in a hybrid approach.
 938 In case of a Stokes-type system, yielding after discretization a saddle point problem, the
 939 strategy can either be applied on the positive definite Schur complement for the pressure or,
 940 as it was done in [139], on the indefinite velocity-pressure system. In that case an all-at-once
 941 multigrid method with an Uzawa-type smoother acting on both solution components turns
 942 out to be most efficient, see [78]. Numerical and algorithmic studies including multiple faults
 943 and large-scale problems with more than $5 \cdot 10^{11}$ degrees of freedom and more than 245000
 944 cores have been demonstrated [138, 139]. The automatic re-coupling strategy is found to be
 945 robust with respect to the fault location and size and also handling multiple fault. In many
 946 scenarios a complete recovery can be achieved with almost no increase in runtime and while
 947 maintaining excellent parallel efficiency.

948 **Adaptive mesh refinement, load balancing, and application level checkpointing**

949 *Adaptive Mesh Refinement* (AMR) functionality and load balancing require similar data
 950 linearization- and transfer functionality as is needed for application level checkpointing. This
 951 is exploited in the waLBerla framework [24, 221, 222] that features an object oriented design
 952 for composing coupled multiphysics simulation software. waLBerla’s load balancing is based
 953 on routines to transfer simulation data between processors so that functionality to serialize,
 954 pack, send, and unpack all relevant data is already available as a by-product of the AMR
 955 functionality. Note that the waLBerla software architecture imposes this structure for Eulerian
 956 mesh based data as well as for Lagrangian particle-based models and it canonically extends
 957 to coupled Eulerian-Lagrangian multiphysics models. For this to work transparently, the
 958 routines for migrating simulation data must be suitably encapsulated. Then this functionality
 959 can be used to write user level checkpoints either on disk or in memory. Note that writing
 960 checkpoints will inevitably imply overheads in memory consumption and communication
 961 time, but that restoring a checkpoint is cheap, since it initially only requires re-activating the
 962 redundantly stored data. This is especially true when in-memory checkpointing is used as
 963 explored and analyzed in [156]. The simple restoration of checkpointed data may of course
 964 lead to load imbalance, but the functionality to redistribute load is also available as part

965 of the parallel AMR functionality. In this sense, user-level checkpointing can be based in a
966 natural, efficient, and straightforward way on the functionality of parallel AMR algorithms
967 combined with load balancing functionality.

968 **Compression techniques to accelerate checkpoint-restart for Krylov-MG solvers**

969 Compressed checkpointing is a possibility to improve the efficiency of classical checkpoint-
970 restart schemes, both in terms of the overhead to generate the checkpoints and to recover
971 the data if an error occurs. The added efficiency mainly comes from a reduced memory
972 footprint which is beneficial for communication and storage. It is particularly efficient if
973 the compression method is tailored to the target application. As an example, in-memory
974 compressed checkpoints combined with LFLR (see Section 3.2.1) for iterative linear solvers,
975 e.g., multigrid preconditioners in Krylov schemes, are described below.

976 *Lossy Compression:* As already mentioned in Section 3.2.2, paragraph ‘Approximated
977 recovery and restart’, initially only the dynamical data, i.e., the approximate solution, are
978 protected. Lossy compression allows a balance between the accuracy of the discretization
979 error of the assembled system and the numerical error within the solver. Specifically in [10],
980 the SZ library [72, 166, 246] is employed, which prefers, by construction, structured data
981 ideally associated with a structured grid. Another important feature is that the compression
982 accuracy can be prescribed and adapted to the situation. Unfortunately, a higher compression
983 accuracy usually leads to a lower compression rate and higher compression time, which is
984 crucial in terms of resilience overhead.

985 Note that multigrid can be interpreted as a lossy compression technique in itself, with a
986 number of mathematical peculiarities that need consideration [115]. Multigrid algorithms use
987 a hierarchy of grids to solve linear systems in an asymptotically optimal way. This hierarchy
988 can be used to restrict, i.e., lossily interpolate, the iterate from fine to coarse grids. Such a
989 lower-resolution representation of the iterate can then be stored as a compressed checkpoint.
990 Conversely, the multigrid prolongation (coarse-to-fine grid interpolation) operator is used to
991 decompress the data. With only small additional computations, the multigrid hierarchy can
992 also be used for error detection.

993 *Recovery:* Several recovery techniques can be devised [10]. As a baseline approach
994 checkpoint-restart is mimicked and the global iterate is simply replaced with its decompressed
995 representation, independently of the compression strategy. The second proposed approach
996 follows the LFLR strategy and re-initializes only the local data that is lost on faulty computing
997 nodes by using checkpoint data stored on neighbouring computing nodes. Contrary to the
998 first approach, this is mostly local and only needs minimal communication to receive a
999 remotely stored backup. In particular, the recovery procedure itself does not involve the
1000 participation of other processes except those sending the checkpointed data. As a worst-case
1001 fallback when the backup data is not sufficient, a third recovery approach is established, which
1002 is still mostly local. Here, an auxiliary problem is solved iteratively with boundary data from
1003 the neighbouring computing nodes. This is similar to the adaptive recovery techniques for
1004 extreme scale multigrid from above or the approximated recovery and restart of Section 3.2.1.
1005 An auxiliary problem is constructed, either by domain decomposition overlap or the operator
1006 structure, and solved with an initial guess based on the checkpoint data to accelerate the
1007 iterative recovery phase. Experiments show that this approach can almost always restore the
1008 convergence speed of the fault-free scenario independently of the used backup technique, only
1009 the number of additional local recovery iterations varies. For more details, we refer to [10].

1010 Resilience with sparse grids

1011 Resilience can be added on various abstraction levels of the algorithm. For PDE problems
1012 one traditionally adds resilience on the level of linear algebra operations, on the solver level
1013 for linear/non-linear equations, or on the time-stepping algorithm. However, this may in
1014 some cases not be coarse-grained enough to minimize the overhead of resilience techniques,
1015 especially when errors occur rarely. In [123, 129, 130, 134, 192, 199] the authors demonstrate
1016 a fault-tolerant framework for solving high-dimensional PDEs that applies fault tolerance
1017 on top of the individual PDE solver. The framework boosts the scalability of black-box
1018 PDE solvers while making it simultaneously resilient to faults by applying the sparse grid
1019 combination technique. In this technique the PDE simulation is distributed over many coarse
1020 grids, which can be processed in parallel. At regular intervals the results of these grids are
1021 combined to obtain the final sparse grid result. In presence of faults the affected grids can be
1022 neglected and an alternative combination scheme is calculated via an optimization routine.
1023 If too many grids are lost, the last combination result serves as an in-memory checkpoint
1024 to recompute the required grids. In [192] it is shown that this lossy recovery provides very
1025 good results even with high error frequencies. At the same time the parallel efficiency is only
1026 slightly affected.

1027 Adaptive mesh refinement

1028 Adaptive refinement techniques in combination with finite element methods are well estab-
1029 lished for fault-free computations. In terms of fault tolerance, this means that in addition to
1030 the assembled linear system, the geometric mesh structure must be protected. This requires
1031 the reconstruction of the data structures containing the mesh hierarchy. For the use of
1032 multigrid or multilevel methods, we also need to recover multiple levels of adaptive grid
1033 refinement after a fault has occurred. The recovery process must take into account the
1034 intra-grid as well as the inter-grid data dependencies.

1035 We refer to [233] for a parallel adaptive multigrid method that uses a sophisticated
1036 dynamic data structures to store a nested sequence of meshes and the iterative evolving
1037 solution. Stals demonstrates that it is possible to implement a fault recovery procedure
1038 that builds on the original parallel adaptive multigrid refinement algorithm [232] in the case
1039 of a fail-stop fault. It is assumed that a copy of the coarsest grid can always be accessed
1040 after a fault has occurred, i.e., it is stored off the processor. The challenge in recovering
1041 an adaptively refined grid is that the mesh distribution changes during any load balancing
1042 procedures, i.e., the local information that was available during the original refinement process
1043 will have been modified or removed. Nevertheless it is demonstrated that the neighbouring
1044 healthy processors contain enough intact information so that the necessary structure can be
1045 recovered to pass into the refinement routine. In the case of uniform refinement, the original
1046 multilevel grid is recovered. In the case of an adaptively refined grid, enough of the structure
1047 is recovered to re-establish the correct communication pattern allowing the solution process
1048 to run to completion, but potentially with reduced accuracy. The neighbouring healthy
1049 processors will only contain enough information to guide the refinement around the edge of
1050 the recovered subgrid. Further refinement within the interior of the recovered subgrid may
1051 be required to improve the accuracy of the solution.

1052 These techniques were implemented with minimal disruption to the original code. An
1053 example of one the few necessary modifications is that in the original code, communication
1054 was used to ensure that the elements were refined in the appropriate order to avoid degenerate
1055 grids. In the resilient version of the the code that communication had to be removed as the

1056 refinement was restricted to the faulty processor.

1057 3.3 Error oblivious algorithms

1058 In this section, we give examples of algorithms that are error oblivious in the sense that they
1059 can recover without assistance from errors that do not occur too frequently. For example,
1060 many fixed point iterative solvers are able to execute to completion if, e.g., a bit flip error
1061 occurs in the solution vector. However, every error likely increases the execution time of
1062 the algorithm. We thus define two quality criteria for error oblivious algorithms and use to
1063 assess the examples in the remainder of this section: (i) correctness, and (ii) efficiency in
1064 terms of execution time.

1065 Finding an algorithm that fulfills (i) and can also compete against error aware algorithms
1066 as described in Section 3.2 remains an open problem.

1067 Error oblivious usually means that an error slowly ‘leaves the system’ during several
1068 iterative sweeps over the data. Error mitigation in error aware algorithms, on the other hand,
1069 requires specific measures to correct the error, and can only be applied when the error has
1070 been detected on a hardware, middleware or algorithmic layer, but removes the disturbance
1071 of the calculation process by the error immediately.

1072 We do not expect the error oblivious algorithms to be impervious to all types of errors.
1073 An iterative method may be not error oblivious if the error changed the matrix entries. This
1074 concept is defined as selective reliability, see Section 3.2.1.

1075 3.3.1 Gossip based methods

1076 A potentially interesting alternative in large-scale parallel environments that does not require
1077 any explicit error detection mechanisms utilizes gossip-based methods and their inherent
1078 resilience properties. Such algorithms by nature build up redundancy in the system and can
1079 thus can efficiently recover automatically from various types of faults/errors without any
1080 need to explicitly detect them. In particular, Gansterer et al. have studied and extended the
1081 resilience of gossip-based aggregation and reduction methods [56, 101, 190]. Based on these
1082 building blocks, they have developed and analyzed several more complex resilient numerical
1083 algorithms, such as orthogonalization methods [101, 102], eigensolvers [235], and least squares
1084 solvers [205].

1085 While the strong resilience properties and execution-time robustness of these approaches
1086 are promising, there is a certain price in terms of basic runtime compared to classical
1087 deterministic numerical high performance algorithms. It remains to be investigated whether
1088 they can be competitive in a fault-prone, but otherwise classical system with global view
1089 and centralized control. Their competitiveness can be expected to increase significantly if
1090 some of these classical properties have to be weakened at the extreme scale.

1091 3.3.2 Fixed-point methods

1092 We view fixed-point methods as methods that converge globally when certain conditions are
1093 satisfied. For example, the Jacobi iterative schemes will converge for any initial guess if the
1094 matrix is diagonally dominant. Fixed-point based iterative methods are by design resilient
1095 to bit flips. However, the convergence delay can be significant. Anzt et al. [12, 13] propose
1096 techniques improving the cost-robustness with little overhead.

1097 A class of numerical algorithms that by design have properties attractive for resilience
1098 are asynchronous iterative methods [23, 29, 39, 40, 97, 230, 231, 242]. In order to avoid

1099 misunderstandings, we point out that this class of methods is unrelated to the idea of
1100 asynchronous dynamic load balancing [155] as addressed in Section 2.1. Instead, asynchronous
1101 iterative methods, stemming from the concept of chaotic iterations [60], are fixed-point
1102 methods that seek the solution of a problem by independently updating subdomains – which
1103 can be subdomains in the geometric sense, subsets, or individual components of the solution
1104 approximation – according to some fixed-point linear or nonlinear iterative scheme. A
1105 particularity of the asynchronous methods is that the independent updates neither adhere to
1106 a specific update order, nor synchronize in terms of a handshake with other updates, but
1107 still converge globally in the asymptotic sense. In particular, these methods are robust with
1108 respect to some subdomains being updated at a much lower pace as each update just uses
1109 the most recent non-local information available. In that sense, asynchronous solvers can have
1110 good performance in unreliable environments where messages can be dropped or processes
1111 can become unresponsive for limited time. Also, in cases where messages are corrupted (and
1112 corruption can be detected), an asynchronous solver can simply drop such a message. In cases
1113 where processes remain unresponsive, a mechanism is still needed to recover that process
1114 and its state, but the remaining processes can continue computing unchanged. Therefore,
1115 asynchronous methods are somehow error oblivious.

1116 With the increasing cost of global synchronizations, and the attractive properties con-
1117 cerning fine-grained parallelization and resilience against communication delays and errors,
1118 asynchronous methods have gained attention in particular for numerical computations [243].
1119 Chow et al. [65, 66] developed an asynchronous algorithm for generating incomplete fac-
1120 torizations, Coleman et al. [68] further improved this algorithm by employing measures
1121 that reduce the runtime overhead when encountering errors. More general is the idea of
1122 asynchronously updating subdomains in Schwarz decompositions. In particular asynchronous
1123 restricted additive Schwarz methods and asynchronous optimized Schwarz methods have been
1124 identified to combine algorithm-inherent resilience with scalability on pre-exascale hardware
1125 architectures [83, 103, 112, 175, 270].

1126 Independently, asynchronous multilevel methods have been proposed and analyzed under
1127 the name Fully Adaptive Multigrid method [214]. Here the multigrid smoothing process is
1128 executed asynchronously so that it can be employed for concurrent operations on different
1129 levels of the mesh hierarchy. The iteration is executed in a Southwell style [228] and
1130 is controlled by efficient hierarchical error estimators [216]. The parallel implementation
1131 [215] will automatically correct errors. More recently, asynchronous methods have been
1132 proposed for nonlinear multi-splitting [244] and eigenvalue computations like Google’s
1133 Pagerank algorithm [157]. More recently, also the idea of asynchronously solving coarse-
1134 grid error correction equations has been investigated, leading to an asynchronous multigrid
1135 algorithm [266]. While case studies reveal attractive properties, these newly developed
1136 asynchronous iterative methods (such as asynchronous multigrid) are not fixed-point iterations,
1137 and developing a convergence theory for those algorithms remains a challenge.

1138 3.3.3 Krylov subspace solvers

1139 A comprehensive overview about the use of selective reliability with Krylov methods in the
1140 presence of bit flips is given in James Elliott’s PhD thesis [86]. Elliott evaluates the CG
1141 and GMRES solvers with the algebraic multigrid preconditioner, see also [10] for a more
1142 recent study. Coleman et al. [68] consider Krylov subspace solvers in combination with the
1143 incomplete ILU algorithm ParILUT. In [84] Elliot et al. investigate the effect of bit flips on
1144 the convergence of GMRES and propose strategies for minimizing the numerical impact.

1145 The authors of [31] present a monotonicity-based fault detection and correction procedure

1146 for a Generalized Conjugate Gradient Krylov solve and perform tests with manual fault
1147 injection. While the solver manages to converge even with large amounts of corrupted data,
1148 the basic recovery procedure speeds up convergence with minimal detection and correction
1149 overhead.

1150 In [218] the authors use a slightly different terminology and call their method numerically
1151 self-stabilizing, a term which originates in the context of distributed systems [77]. They
1152 introduce two error oblivious [77] iterative linear solvers: one for the steepest descent and one
1153 for conjugate gradient. In the latter case, they consider necessary conditions for conjugate
1154 gradient to converge. Those conditions are borrowed from non-linear conjugate gradient [276]
1155 and are maintained in a correction step (typically performed every other ten iterations). The
1156 correction step does not explicitly correct errors, but re-computes quantities such as the
1157 residual at regular intervals. Therefore, we classify these methods as error oblivious instead
1158 of error aware.

1159 **3.3.4 Domain decomposition**

1160 In [116] Griebel and Oswald use probabilistic analysis to model the effect of errors on the
1161 convergence of the classical overlapping Schwarz algorithm. They conclude that this method
1162 does indeed converge in the presence of errors. Glusa et al. [113] mention that asynchronous
1163 domain decomposition methods are by definition fault-tolerant. In [184,210,211], the authors
1164 discuss resiliency of a task-based domain decomposition preconditioner for elliptic PDEs. By
1165 leveraging the domain decomposition approach, the problem is reformulated as a sampling
1166 problem, followed by a regression-based solution update. The regression is formulated and
1167 implemented such that it is oblivious to corrupted samples. The authors combine this
1168 algorithmic approach with a server-client implementation based on ULFM, see Section 2.2.2.
1169 They show promising results of this approach in terms of resiliency to missing tasks, corrupted
1170 data and hardware failure.

1171 **3.3.5 Time stepping**

1172 In [119], iterative time-stepping using spectral deferred corrections are shown to be error
1173 oblivious at the cost of more iterations for the affected time-step. With error-estimators in
1174 place, time-integration techniques like Runge-Kutta methods will repeat the calculation of a
1175 time-step with smaller step sizes, if errors in the solution vectors are relevant [61]. This type
1176 of algorithms is resilient against errors in the solution vector of the new time step. Repeating
1177 the new time step with a reduced time step size is not the optimal measure in case of an
1178 error where repeating the step with the same time step size would be more efficient, but it
1179 leads to correct results.

1180 **4 Future directions**

1181 In the final section we focus on the future direction of resilient algorithms. We highlight
1182 what changes need to be made to current infrastructures to support the goals proposed by
1183 algorithm and application developers. Furthermore, we list those algorithms that are likely
1184 to come to the forefront as resiliency plays a more important role in the cost-benefit analysis
1185 of extreme scale simulations. And we mention some numerical methods that are yet to be
1186 fully explored in the context of resilient algorithms.

1187 4.1 Systems in support of resilient algorithms

1188 We propose that resiliency will only be obtained by a multilayered approach incorporating
1189 operating systems, file systems, communication, programming models, algorithms, applica-
1190 tions and education. In terms of the layers covered by infrastructure, the goal is to increase
1191 systems and delivered performance while keeping the detectable errors in the upper algorithm
1192 based layers constant. We refer the reader to the recently published report by Radojkovic et
1193 al. [207] for an overview of the needs of the next generation HPC systems.

1194 4.1.1 Error correcting codes

1195 Poulos et al. [204] propose hardware ECC assistance that can pass error syndrome information
1196 through to an application and use this to fix detected errors. When an ECC hardware
1197 error occurs that results in a *Detectable, but Uncorrectable Error* (DUE), the ECC hardware
1198 generates a syndrome which is a byproduct of the error detection. For many ECC schemes,
1199 a syndrome that corresponds to a DUE can be used to generate a list of possible corrections,
1200 one of which is taken to be the original uncorrupted data. In this work, the authors show
1201 that this set is relatively small, meaning that the set of potential values for an application to
1202 search for their correct answer (before corruption) is also small. They also study the error
1203 value distribution and show that for certain classes of problems it can be easy to identify
1204 obviously wrong answers. For the application studied in [204], work was done to correct a
1205 hydrodynamics application using conservation laws and average of neighbor cells. This work
1206 requires changes to the hardware error reporting techniques and modification to the operating
1207 system to determine which application observed the DUE and pass it to an interrupt handler.

1208 4.1.2 Improving checkpoint/restart

1209 Independent of any additions, changes or new developments in the algorithmic or the
1210 system area, checkpoint/restart will remain a necessary component for any system. For one,
1211 no other technique can provide the needed resilience against full system outages; further,
1212 checkpoint/restart is also needed for developers to deal with limited job execution times and
1213 possible migration between systems or debugging purposes at large scale.

1214 Improving classical checkpoint/restart for homogeneous systems

1215 Observing the necessity of checkpoint/restart makes it critical to further optimize, enhance and
1216 support efficient checkpoint/restart mechanisms—even on classical, homogeneous systems—
1217 and provide users with library based solutions for core checkpoint/restart functionality. In
1218 particular, the following avenues should be pursued to optimize checkpoint/restart.

- 1219 ■ Use additional algorithmic techniques to be able to reduce checkpoint frequency.
- 1220 ■ Reduce data to be written to disk by eliminating redundancy and possibly compressing
1221 checkpoint information. Note that suitable data compression will typically require
1222 user-level knowledge, suitable interfaces must be provided.
- 1223 ■ Overlap/Offload checkpoint operations to allow for asynchronous checkpoint/restart
1224 operations.
- 1225 ■ Integrate checkpoint/restart with novel programming approaches to minimize check-
1226 pointable state.
- 1227 ■ Keep the restart requirements local to the neighbour nodes of the failed node.

XX:32 Acronyms

- 1228 ■ Localize checkpoint data to own or localized nodes. This could be supported by local
1229 non-volatile memory, as targets for checkpoint data. While this has the potential to
1230 reduce communication, as it avoids remote data transfers, it may require additional
1231 hardware support to retrieve data from non-functional nodes, e.g., by accessing data
1232 through fast JTAG-like interfaces.
- 1233 ■ In memory checkpointing.
- 1234 ■ Exploit user-level knowledge for serializing, packing, compressing data, see e.g. how exist-
1235 ing AMR functionality [156] can be exploited for efficient checkpointing in Section. 2.2.1.

1236 Checkpoint/Restart for heterogeneous systems

1237 In addition to classical checkpoint/restart for homogeneous systems, node-local check-
1238 point/restart support for heterogeneous systems will help containing error and failure
1239 propagation. Such support may be provided transparently to the application by the under-
1240 lying infrastructure, such as GPU drivers or task-based environments, or exposed in the
1241 programming model, such as OpenMP Offload [76].

1242 4.1.3 Scheduler and resource management

1243 Support for resilience, especially at the workflow-level, has a direct impact on resource
1244 management in HPC systems and hence requires new developments in this area as well.

1245 Node-level parallelism

1246 With increasing node-level parallelism, the impact of OS noise (typically caused by un-
1247 predictable interrupts) becomes even more important. Therefore, dedicated node-level
1248 resources are needed to exclusively run the OS and minimize the impact of OS noise on the
1249 multi-threaded application running on the other cores.

1250 Adaptive system and application load balancing

1251 The batch scheduler needs to adaptively balance the system load onto the available resources,
1252 via seamless application migration. While the application needs to adapt to the capabilities
1253 of the newly allocated resources, if different from the original allocation, without incurring
1254 performance penalties. The former has typically been implemented via checkpointing and
1255 process migration [174]. The latter has typically been implemented for applications that can
1256 adjust their granularity, e.g. from finer to coarser, depending on resource availability either
1257 triggered by the application or the system [46]. When exposing and expressing parallelism
1258 in applications, in addition to accounting for and matching the multiple levels of hardware
1259 parallelism (nodes, sockets, cores), the decomposition granularity needs to be flexible to
1260 support evolvability and malleability and allow for adaptive load balancing at the application
1261 and system levels.

1262 Adaptive resource management

1263 The batch scheduler in conjunction with the distributed runtime system employed by the
1264 application (e.g., MPI, Charm++, HPX) needs to support resources errors/failures and
1265 recover them without terminating the applications in the process. This approach should work
1266 both with rigid and moldable applications as well as with evolving and malleable applications.

■ **Table 2** Properties of numerical algorithms fostering or helping resilience

categories	solvers	discretization
redundancy	×	×
replication	×	
hierarchical methods	×	×
mixed precision	×	×
error control	×	×
locality-emphasizing schemes		×
asynchronous methods	×	×
embarassingly parallel	×	
stochastic > deterministic		×
iterative vs direct solvers	×	
matrix-free / low memory footprint	×	×

1267 4.2 Programming models with inherent resiliency support

1268 Certain applications and algorithms may naturally be resilient against errors. This may lend
 1269 them as natural candidates for asynchronous parallel execution (via asynchronous many-task
 1270 programming). While this mitigates the challenges associated with bulk synchronous parallel
 1271 execution, asynchronous parallel execution may influence, in the presence of silent errors, the
 1272 convergence rate of the numerical algorithms and might lead to incorrect results.

1273 Programming model and runtime support for resilience can offer transparent handling of
 1274 errors and failures or can assist the application in handling them. Consistent programming
 1275 model support for resilience based on realistic error/failure models is needed to properly
 1276 handle such events with low overhead. Higher-level abstractions for programming resilient
 1277 applications are needed to help with error/failure handling complexities and to offer reuse of
 1278 concepts and codes across applications.

1279 4.3 Future directions for the solution of partial differential equations

1280 In this section, we focus on discretizations for linear and non-linear partial differential
 1281 equations as well as solvers for the resulting discrete and sparse systems of equations. We
 1282 introduce a list of algorithmic properties that we found are, or can be, contributing to
 1283 the resilience of the algorithms described in Section 3. Table 2 lists these properties and
 1284 indicates where we found relevant examples of how they can foster resilience for either linear
 1285 or non-linear solvers or for spatial or time discretization. In the following subsections we
 1286 describe these examples in more detail and highlight the several (mutually related) properties
 1287 that could be of interest in the context of resilient algorithms.

1288 4.3.1 Redundancy and replication

1289 A failure that is not fixed by the system (hardware and middleware) typically results
 1290 in a loss or corruption of data. To tackle this problem, redundancy techniques can be
 1291 used to detect and recover from data corruption and data loss. The performance of these
 1292 algorithms is usually measured in the amount of memory and computational overhead
 1293 they entail, the detection rate of errors, the rate of false-positives they achieve, and the
 1294 accuracy of the recovery. Optimizing these performance indicators should be of main concern
 1295 for future algorithm design. One existing class of algorithms that apply redundancy are

1296 multiresolutional techniques such as multigrid and the sparse grid combination technique
1297 described in Section 3.2. They inherently add redundancy through the hierarchical structure.
1298 Sparse grid combination techniques calculate the same solution on different anisotropic grids.
1299 The coefficients of the combinations of the components grids can be recalculated if one
1300 or more nodes are lost due to faults. This redundancy of the component grids allows the
1301 algorithm to obtain an alternative approximation of the solution. However, if a component
1302 grid is distributed on too many nodes, then the approximation will fail if a fault occurs on
1303 any one of those nodes. Another class of algorithms add redundancy through recomputation
1304 with different models and configurations such as in ensemble or multifidelity techniques. A
1305 more straight-forward approach is to directly add redundancy through replication of certain
1306 algorithmic paths, cf. the following subsection on recalculation techniques.

1307 Depending on the underlying architecture, replication can be a competitive option to
1308 increase detected and undetected error robustness. If computation speed significantly outpaces
1309 memory access and communication, each operation can be executed multiple times while the
1310 data is still accessible in the RAM. This can be used for redundancy-based sanity checks of
1311 low-level operations or even for checksum-like approaches.

1312 Overlapping data in parallel algorithms can serve as a starting point for mitigation,
1313 albeit not for detection. In the case studies explored in Section 3.2.2, these are applied to
1314 elliptic PDEs, though an extension to other models should be feasible. Furthermore by even
1315 increasing the ghost layer size and thereby adding extra redundancy, other reconstruction
1316 possibilities might become possible. This could already be taken into account during the
1317 domain partitioning process.

1318 4.3.2 Hierarchy and mixed precision

1319 Hierarchical discretizations have proven to be advantageous in various respects. Related
1320 notions are multi-resolution or multi-level discretizations, but also (recursive) sub-structuring
1321 in the engineering nomenclature of the *Finite Element Method* (FEM). Built into the hierarchy
1322 are problem-inherent information and structures that are well-suited for modern hierarchy-
1323 based solvers. In FEM, for example, hierarchical bases carry information about both location
1324 and frequency, which leads to a special built-in redundancy that can be exploited for error
1325 detection (see Section 3.1). Therefore, from a resilience perspective, hierarchy should be a
1326 core paradigm for discretization design. This applies irrespectively of whether the hierarchical
1327 bases are formulated in the spatial (h) or the order (p) sense.

1328 From a solver perspective, multigrid methods for elliptic and parabolic PDE problems are
1329 relevant approaches towards resilient numerical algorithms. They inherently act on different
1330 granularities, representations, scales, and levels and can be used to quantify differences
1331 between these levels. For local recovery, local multigrid methods are highly efficient, especially
1332 when they can be accelerated with the superman strategy [138]. Additionally the low-
1333 resolution duplicates can be used for some kind of approximate recovery or minimal rollback
1334 like re-application of the smoother on a specific level in a multigrid scheme. Detection
1335 of errors within multigrid is often possible due to algebraic relations or on the basis of
1336 hierarchical multi-grid-inherent error estimates [11, 139, 216], which hold true inside such
1337 schemes. As stated in Section 4.3.1, the inherent redundancy incorporated in these algorithms
1338 is also beneficial.

1339 Mixed-precision arithmetics are typically used within the numerical solver parts to speed-
1340 up computations. However, the discretization can enable the flexibility to store data at
1341 varying precision. Examples for this are hierarchical approaches such as hierarchical bases,
1342 where a function value is stored as a hierarchical surplus only. As another example, the

1343 usage of wavelets in multiresolutional analysis can serve. In both cases, contributions of
1344 higher levels typically require less accuracy, as only the most significant bits contribute to
1345 the overall point values.

1346 4.3.3 Error control

1347 For many numerical methods, a wide range of classical a priori and a posteriori error
1348 estimation techniques are available, see among many others [6, 22, 122, 152, 160, 206, 216],
1349 which constitute the basis of many adaptive numerical algorithms.

1350 Adaptive time discretization methods are the state of the art for ODE solvers, while, for
1351 PDE solvers, spatial adaptivity techniques are also widely used. Local time step adaptation
1352 is feasible in the framework of so called local time stepping or multirate approaches, where
1353 different components of the system can have different time step sizes, see [43, 53, 94, 106, 217,
1354 220, 224], which are however still far from mainstream for most applications. For PDE solvers,
1355 local spatial adaptivity techniques are also very common [20, 21], but their incorporation in
1356 operational applications is still a research topic, see e.g. [36, 163, 185, 203, 255] for developments
1357 concerning oceanography and numerical weather forecasting.

1358 The error estimations on which all these methods rely on also constitute the basis of an
1359 error detection mechanism, since some undetected errors, like bit flips on significant floating
1360 point digits, will result in errors exceeding the allowed error tolerances. To some extent,
1361 these techniques are also examples of ABFT or error oblivious approaches, since bit flips and
1362 other silent errors occurring during the computation of the solution at the next time step or
1363 on a refined mesh could be automatically corrected by the repeated computations triggered
1364 by the error threshold violation. Furthermore, silent errors in the data at the current time or
1365 mesh level could be identified by the failure of the time step or mesh refinement to correct
1366 the error.

1367 Combined with other ABFT strategies, adaptive discretization strategies based on error
1368 estimators can be a powerful and so far rather underrated tool for protecting a simulation
1369 from undetected errors in the solution vectors. On the other hand, error estimators should not
1370 be used as a black box for resiliency purposes. Indeed, errors can lead to severe over-resolution
1371 or, potentially, even under-resolution in space or time and the error estimators themselves
1372 could be affected by undetected errors.

1373 As seen in Section 3.1.4, some iterative solvers for the solution of linear systems have
1374 invariants, such as monotonicity for Krylov solvers. These properties can be put to good use
1375 in devising resilience strategies, for example activating an additional restart of the Arnoldi
1376 procedure as soon as an increase in the residual norm is observed.

1377 The idea of interval arithmetic is to compute bounds of intervals that always contain the
1378 exact result [7, 158]. Probabilistic methods for rounding error estimation [71, 96, 98, 198, 258]
1379 require several executions of arithmetic operations with different perturbations or different
1380 rounding modes (for instance three executions for Discrete Stochastic Arithmetic [81]). With
1381 both approaches, the comparison of several computed results enables one to control rounding
1382 errors (or detect and mitigate actually wrong results).

1383 4.3.4 Locality, asynchronicity and embarrassingly parallelism

1384 One important aspect of resilient algorithms is error confinement as global dependencies
1385 propagate errors to other processors and complicate recovery. Locality-emphasizing numerical
1386 algorithms achieve this by limiting dependencies to local areas or completely removing them.
1387 Consequently, error mitigation can be limited to a local subdomain. Typical examples for

1388 these schemes are domain decomposition, which splits the domain into several subareas, and
1389 classical discretization schemes such as finite elements, finite differences and finite volumes.

1390 As mentioned in Section 3.2.1, domain decomposition schemes such as additive Schwarz
1391 methods, or substructuring-inspired FETI [90] or also the fully adaptive multigrid method
1392 [214] are naturally asynchronous and resilient to message loss. In this context, we use the
1393 term asynchronous primarily in the sense of reducing the time synchronicity in parallel
1394 computations – from communication-avoiding schemes via a reduction of synchronization
1395 points up to vastly decoupled schemes. Using this inherent property, a failure in a subdomain
1396 would result in a message loss that does not hinder convergence in other subdomains, because
1397 a global wait for a message update and synchronization are not necessary. In addition,
1398 asynchronous methods may better adapt to heterogeneous processors and networks than
1399 their synchronous counterparts as it has been shown in the context of Grid computing [19,59].
1400 Both the localized and asynchronous approaches, achieve their impact through a decoupling of
1401 computations. Going further in this direction leads to embarrassingly or nearly embarrassingly
1402 parallel algorithms. These represent a group of algorithms where it is relatively easy to
1403 decouple subproblems in time or space. The subproblems can therefore be calculated
1404 completely independently, and errors do not propagate to other subproblems. Examples of
1405 such methods are Monte Carlo simulations and computations with the sparse grid combination
1406 technique. Since it is expected that only a few tasks will encounter errors and the scheduling
1407 is automatically balancing the load, the overall execution time does not suffer too much.
1408 Future algorithmic design should therefore aim at increasing asynchronicity and locality to
1409 move towards embarrassingly parallel problems.

1410 **4.3.5 Stochastic**

1411 Stochastic methods can be superior to deterministic methods when it comes to resilience.
1412 Stochastic methods do not require the program to take a deterministic path, faulty parts
1413 can be neglected or exchanged easily by other results. A popular example are Monte Carlo
1414 methods where we sample randomly in the computation domain and can simply neglect
1415 failed samples. Ensemble methods are examples where different instances or models of a
1416 concrete problem setting are computed. Even if one of these computation fails, the ensemble
1417 computation can still return a – maybe slightly less accurate – result. Stochastic elements
1418 can therefore help the future algorithm design to reduce the dependencies on specific results
1419 of the computation. These methods, however, need to be evaluated not just by highlighting
1420 their resilience properties, but also taking into account the cost of a single run: if a single
1421 run is expensive to complete, simply discarding it might be impractical.

1422 **4.3.6 Iterative methods**

1423 Iterative solvers may be viewed as inherently more robust than direct solvers because they
1424 do not compute their solution using a pre-defined sequence of numerical operations as direct
1425 solvers typically do. Indeed, by their nature, they perform a sequence of operations to update
1426 and improve their current approximation. If an error is encountered during computation, the
1427 probability of deleting this error or at least its effect may be higher than in a direct solver.
1428 Especially fixed-point-based methods (domain decomposition, relaxation, ...) may be viewed
1429 as inherently resilient as they have the property to always converge to the correct solution
1430 independent of the initial state (global convergence). Some errors may induce a low influence
1431 on convergence speed and can thus be safely ignored. In other cases, a restart – optionally
1432 with recovery techniques – may be employed to ensure both resilience and efficiency in terms

1433 of runtime.

1434 **4.3.7 Low memory footprint – matrix-free**

1435 The classical approach to represent linear operators as sparse matrices produces large
1436 amounts of static data which has to be restored upon failure. Checkpoint-restart approaches
1437 feature high memory cost, naturally multiples of the storage needed for the solution vector.
1438 Algorithmic alternatives to checkpoint-restart require possibly complicated or costly re-
1439 assembly. Matrix-free methods do not represent the operators as static data in the first place.
1440 Therefore, large sparse matrix data structures do not have to be restored upon failure as they
1441 are computed on the fly anyway. Extreme-scale applications will benefit from matrix-free
1442 approaches due to their low memory footprint, also in terms of runtime, (due to high memory
1443 access cost) and higher limits for the overall problem size [25, 27].

1444 In addition to saving memory and, therewith, reducing the risk of memory corruption,
1445 matrix-free methods can also be combined with automatic code generation [162] in a stencil-
1446 based approach, i.e., for finite difference methods on uniform structured grids. In such cases,
1447 the matrix entries may be ‘hard wired’ into code, such as 5-point stencils for Laplace’s equation.
1448 Automatic code generation provides a means to increase resiliency in the code generator or
1449 domain specific language and, thus, facilitate resilience aware software development.

1450 For finite element methods, one can use local assembly kernels [26]. Here, the trade-off
1451 between computation and storage and, in the future, resilience is relevant in particular for
1452 higher order elements.

1453 **4.4 The final mile: towards a resilient ecosystem**

1454 The future directions described above will provide critical enhancements towards providing
1455 resilient computation for numerical simulations. Alone, however, they are insufficient, as
1456 they must be embedded in the larger ecosystem and in the efforts to make that ecosystem
1457 support such novel resilience approaches. This requires another set of crucial developments.

1458 **4.4.1 Tools to support resilience software development**

1459 Developers will need the right tools to support their algorithmic efforts. These tools, as
1460 they exist today, are often designed without faults and errors in mind and, therefore, do
1461 not sufficiently support the development of resilient systems. In particular, we identified
1462 three areas in which enhanced tool support for resiliency is needed: a) introspection to help
1463 track errors and failures along with their root causes, b) validation through controlled fault
1464 scenarios to enable targeted testing of new error mitigation features, and c) transformation
1465 to transparently add error and failure checks into codes.

1466 **Tools for introspection**

1467 Introspection is critical to ensuring early error detection and the timely activation of correction
1468 and mitigation mechanisms throughout the various layers of the software ecosystem.

1469 *System Monitoring:* Knowing about the health state of a system requires monitoring it
1470 and understanding its behavior. Future work needs to focus on scalable system monitoring,
1471 real-time analyzes of system monitoring data, and autonomous decision making on corrective
1472 actions for self-aware resilient systems. In order to gain a deeper understanding, types of
1473 monitored data should be homogenized across system and sites, and, if possible, sanitized
1474 logs should be available to the community.

1475 *Application and Data Structures Monitoring:* Applications need to automatically monitor
1476 their performance and correctness with the use of tools. The tools can be developed in
1477 abstraction, at the compiler-level, or at the runtime-level.

1478 **Tools for validation**

1479 Currently, there are no standard tools to test the correctness and performance of resilient
1480 algorithms under undetected errors and fault. This is due to a lack of fault injection tools
1481 that reflect realistic situations. DeBardeleben et al. [120] have developed a hardware error
1482 simulator tool to understand the behavior of numerical algorithms under faulty hardware with
1483 a great accuracy, but this approach cannot evaluate the execution time of resilient algorithms
1484 at scale. Vendors provide fault injection tools [126, 148] for better execution efficiency,
1485 compromising the accuracy of the hardware behavior. Compiler approaches or other in-house
1486 error injections [49, 107] could allow the program to execute as efficiently as the original
1487 binary, but the correctness is further compromised. There are also tools that can analyze
1488 an application's vulnerability very quickly but do not actually produce the application's
1489 faulty output. One technique for this, DiSCvar [177], uses algorithmic differentiation and
1490 exposes how changes to each variable impact output results. It is important to note that
1491 these techniques do not actually produce that corrupted output. Hence, they are very fast
1492 but they may not be useful to developers looking to explore precisely how corruption changes
1493 their application. It is likely that a combination of these techniques, which identify most
1494 critical regions of an application coupled with fault injection at those locations, may serve as
1495 a good compromise between the two techniques.

1496 Any novel approaches that fill the gap between the accuracy and execution efficiency of
1497 error injections will facilitate the code development of resilient algorithms, and the new tools
1498 should be built with the existing continuous integration infrastructure. Such tools likely
1499 require hardware knowledge that is considered intellectual property by the semiconductor
1500 vendors. However, efforts which explore this space using open hardware technologies (RISC-V,
1501 Sparc, etc.) can shed light on this space but may be of varying usefulness when application
1502 developers look to understand how their applications will perform on hardware that has not
1503 been fault injected at the register transfer or microcode level.

1504 **Tools for code transformation**

1505 Compilers are able to generate binaries with resilience capability as suggested in the work
1506 by [209]; the generated binary instruments redundant computation, register allocations
1507 to enable error detection and correction during program execution. The recent work by
1508 Lin [170] leverages LLVM to generate SIMD instructions to perform redundant computation
1509 and verification. Source-to-source code transformation has been proposed to enable triple
1510 modular redundancy in loops [168] and automatic instrumentation of checkpointing [212].
1511 Similarly, this idea can be extended to redundant threading for error mitigation, facilitated
1512 with OpenMP-like programming language extension [140]. These approaches automatically
1513 introduce resilience with some performance penalty, preventing the users from selective
1514 adaptation of resilience for performance optimization, and these redundant computations are
1515 benefited from the memory hierarchy, preventing doubling (or tripling) of the execution time.

1516 In addition to such specific systems that support the addition of resilience to existing
1517 codes, automated generation of code, e.g., via *Domain Specific Languages* (DSL) can help
1518 with the transparent support of resilient computation. Examples for this can be stencil
1519 generators, as already discussed in Section 4.3.7.

1520 4.4.2 User/Programmer education

1521 According to the system log study by [75], many application job failures are triggered by
1522 the mistakes of the users such as script errors and program bugs including excessive file and
1523 thread creations. This means that better software engineering practices and training of users
1524 should be pursued with similar efforts to the deployment of resilience strategies.

1525 The Exascale Computing Project (ECP) by the US DOE has made a substantial investment
1526 on educating tools, software engineering and HPC system usage for a variety of the users.
1527 Additionally, the scientific and mathematical library teams in the ECP have introduced
1528 software engineering policies [272] to improve the software quality, documentation and testing
1529 process for better interoperability and composability of multiple library packages. This
1530 activity, though not directly relevant to resilience, will gradually help to reduce application
1531 errors and failures for large scale HPC systems.

1532 5 Conclusions

1533 This article presents a snapshot of current research on resilience for extreme scale computing.
1534 It has grown out of the Dagstuhl seminar 20101 held March 1-6, 2020, bringing experts from
1535 the field together on the topic *Resiliency in Numerical Algorithm Design for Extreme Scale*
1536 *Simulations*. This seminar became a starting point to develop a synthesis between the system
1537 perspective on resilience and the algorithmic perspective.

1538 While resilience is undoubtedly an issue for extreme scale computing, it is less clear
1539 what algorithms on the user or application level can contribute to mitigate faults. The
1540 seminar provided ample room to discuss these topics and thus became the starting point
1541 for this article. Many diverse aspects were found to be relevant, that require a holistic and
1542 multidisciplinary approach involving different and complementary scientific communities.

1543 In particular, it clearly appeared that a fundamental distinction lies in whether faults
1544 are detected or not, and if they are not automatically detected, whether they are detectable.
1545 If they are, algorithms can often be developed to detect errors and in a second stage to
1546 correct them. It was found that some algorithms are naturally tolerant against faults or have
1547 the intrinsic feature to be error oblivious. They can thus be naturally applied on a system
1548 subject to errors.

1549 Besides redundancy and checkpointing as classical techniques to mitigate faults, new
1550 algorithm-based resilience techniques have been developed for several classes of numerical
1551 algorithms. This includes linear algebra and solvers for partial differential equations, two
1552 classes of algorithms that are prominent in many scientific workloads on supercomputers.
1553 Some of these mitigation methods show remarkable success in the sense that faults can be
1554 compensated algorithmically by recovery procedures with only little extra cost in time or
1555 in silicon. On the other hand it also becomes clear that integrating such techniques in a
1556 computational infrastructure is still facing many obstacles. This includes the still poorly
1557 defined interface between user-level fault mitigation techniques and system level functionality,
1558 as, it is, e.g., necessary to reliably and quickly detect a device (core, memory, ...) failure on
1559 a large parallel machine.

1560 Despite its breadth, the article is far from being comprehensive. The selection of topics is
1561 a subjective overview of current research in the field of resilience for extreme scale computing
1562 and it delivers an outlook into possible and promising future research topics and solutions.

1563 — References —

- 1564 1 Emmanuel Agullo, Siegfried Cools, Emrullah Fatih-Yetkin, Luc Giraud, Nick Schenkel, and
1565 Wim Vanroose. On soft errors in the Conjugate Gradient method: sensitivity and robust
1566 numerical detection - revised. Research Report RR-9330, Inria, March 2020. URL: <https://hal.inria.fr/hal-02495301>.
1567
- 1568 2 Emmanuel Agullo, Luc Giraud, Abdou Guermouche, Jean Roman, and Mawussi Zounon. Nu-
1569 merical recovery strategies for parallel resilient Krylov linear solvers. *Numerical Linear Algebra*
1570 *with Applications*, 23(5):888–905, August 2016. URL: <https://hal.inria.fr/hal-01323192>,
1571 doi:10.1002/nla.2059.
- 1572 3 Emmanuel Agullo, Luc Giraud, Pablo Salas, and Mawussi Zounon. Interpolation-restart
1573 strategies for resilient eigensolvers. *SIAM Journal on Scientific Computing*, 38(5):C560–C583,
1574 2016. URL: <https://hal.inria.fr/hal-01347793>, doi:10.1137/15M1042115.
- 1575 4 Emmanuel Agullo, Luc Giraud, and Mawussi Zounon. On the resilience of parallel sparse hybrid
1576 solvers. In *HiPC 2015 - IEEE International Conference on High Performance Computing*,
1577 Bangalore, India, December 2015. URL: <https://hal.inria.fr/hal-01256316>.
- 1578 5 Mark Ainsworth and Christian Glusa. Is the multigrid method fault tolerant? the multilevel
1579 case. *SIAM Journal on Scientific Computing*, 39(6):C393–C416, 2017.
- 1580 6 Mark Ainsworth and J Tinsley Oden. *A posteriori error estimation in finite element analysis*,
1581 volume 37. John Wiley & Sons, 2011.
- 1582 7 G. Alefeld and J. Herzberger. *Introduction to interval analysis*. Academic Press, 1983.
- 1583 8 Md Mohsin Ali, Peter E. Strazdins, Brendan Harding, and Markus Hegland. Complex scientific
1584 applications made fault-tolerant with the sparse grid combination technique. *International*
1585 *Journal of High Performance Computing Applications*, 30(3):335–359, 2016.
- 1586 9 Md Mohsin Ali, Peter E. Strazdins, Brendan Harding, Markus Hegland, and Jay W Larson.
1587 A fault-tolerant gyrokinetic plasma application using the sparse grid combination technique.
1588 In *Proceedings of the 2015 International Conference on High Performance Computing &*
1589 *Simulation (HPCS 2015)*, pages 499–507, Amsterdam, The Netherlands, July 2015.
- 1590 10 Mirco Altenbernd, Nils-Arne Dreyer, Christian Engwer, and Dominik Goddeke. Towards
1591 local-failure local-recovery in PDE frameworks. In *Proceedings of HPCSE’19*. Springer, 2020.
1592 accepted.
- 1593 11 Mirco Altenbernd and Dominik Goddeke. Soft fault detection and correction for multigrid. *The*
1594 *International Journal of High Performance Computing Applications*, 32(6):897–912, November
1595 2018. doi:10.1177/1094342016684006.
- 1596 12 Hartwig Anzt, Jack Dongarra, and Enrique Quintana-Ort. Fine-grained bit-flip protec-
1597 tion for relaxation methods. *Journal of Computational Science*, 36:100583, 2019. URL:
1598 <http://www.sciencedirect.com/science/article/pii/S1877750316303891>, doi:[https://](https://doi.org/10.1016/j.jocs.2016.11.013)
1599 doi.org/10.1016/j.jocs.2016.11.013.
- 1600 13 Hartwig Anzt, Jack Dongarra, and Enrique S. Quintana-Ort. Tuning stationary iterative
1601 solvers for fault resilience. In *Proceedings of the 6th Workshop on Latest Advances in Scalable*
1602 *Algorithms for Large-Scale Systems*, New York, NY, USA, 2015. Association for Computing
1603 Machinery. doi:10.1145/2832080.2832081.
- 1604 14 Rizwan Ashraf, Saurabh Hukerikar, and Christian Engelmann. Pattern-based modeling of
1605 multiresilience solutions for high-performance computing. In *Proceedings of the 9th ACM/SPEC*
1606 *International Conference on Performance Engineering (ICPE) 2018*, pages 80–87, Berlin,
1607 Germany, April 9-13, 2018. ACM Press, New York, NY, USA. doi:10.1145/3184407.3184421.
- 1608 15 Rizwan A. Ashraf and Christian Engelmann. Performance efficient multiresilience using
1609 checkpoint recovery in iterative algorithms. In *European Conference on Parallel Processing*,
1610 pages 813–825. Springer, 2018.
- 1611 16 B. Austin, E. Roman, and X. Li. Resilient matrix multiplication of hierarchical semi-separable
1612 matrices. In *Proceedings of the 5th Workshop on Fault Tolerance for HPC at eXtreme Scale*
1613 *(FTXS)*, pages 19–26, Portland, Oregon, June 2015. doi:10.1145/2751504.2751507.

- 1614 17 Todd M. Austin. DIVA: A reliable substrate for deep submicron microarchitecture design. In
1615 *Proceedings of the 32nd Annual ACM/IEEE International Symposium on Microarchitecture,*
1616 *MICRO 32, IEEE Computer Society, Washington, DC, USA, 1999.*
- 1617 18 A. Avizienis, J. . Laprie, B. Randell, and C. Landwehr. Basic concepts and taxonomy of
1618 dependable and secure computing. *IEEE Transactions on Dependable and Secure Computing,*
1619 1(1):11–33, Jan 2004. doi:10.1109/TDSC.2004.2.
- 1620 19 J. M. Bahi, S. Contassot-Vivier, and R. Couturier. Evaluation of the asynchronous iterative
1621 algorithms in the context of distant heterogeneous clusters. *Parallel Computing,* 31(5):439–461,
1622 2005.
- 1623 20 Wolfgang Bangerth, Carsten Burstedde, Timo Heister, and Martin Kronbichler. Algorithms and
1624 data structures for massively parallel generic adaptive finite element codes. *ACM Transactions*
1625 *on Mathematical Software (TOMS),* 38(2):1–28, 2012.
- 1626 21 Wolfgang Bangerth and Rolf Rannacher. *Adaptive finite element methods for differential*
1627 *equations.* Birkhäuser, 2013.
- 1628 22 Randolph E. Bank and R Kent Smith. A posteriori error estimates based on hierarchical bases.
1629 *SIAM Journal on Numerical Analysis,* 30(4):921–935, 1993.
- 1630 23 G. Baudet. Asynchronous iterative methods for multiprocessors. *Journal of the Association*
1631 *for Computing Machinery,* 25:226 – 244, 1978.
- 1632 24 Martin Bauer, Sebastian Eibl, Christian Godenschwager, Nils Kohl, Michael Kuron, Christoph
1633 Rettinger, Florian Schornbaum, Christoph Schwarzmeier, Dominik Thönnies, Harald Köstler,
1634 et al. walberla: A block-structured high-performance framework for multiphysics simulations.
1635 *Computers & Mathematics with Applications,* 2020.
- 1636 25 Simon Bauer, Daniel Drzisga, Marcus Mohr, U Rüe, Christian Waluga, and Barbara Wohlmuth.
1637 A stencil scaling approach for accelerating matrix-free finite element implementations. *SIAM*
1638 *Journal on Scientific Computing,* 40(6):C748–C778, 2018.
- 1639 26 Simon Bauer, Markus Huber, S Ghelichkhan, Marcus Mohr, Ulrich Rüde, and B Wohlmuth.
1640 Large-scale simulation of mantle convection based on a new matrix-free approach. *Journal of*
1641 *Computational Science,* 31:60–76, 2019.
- 1642 27 Simon Bauer, Marcus Mohr, Ulrich Rüde, Jens Weismüller, Markus Wittmann, and Barbara
1643 Wohlmuth. A two-scale approach for efficient on-the-fly operator assembly in massively parallel
1644 high performance multigrid codes. *Applied Numerical Mathematics,* 122:14–38, 2017.
- 1645 28 Leonardo Bautista-Gomez, Seiji Tsuboi, Dimitri Komatitsch, Franck Cappello, Naoya
1646 Maruyama, and Satoshi Matsuoka. FTI: high performance fault tolerance interface for hybrid
1647 systems. In *Proceedings of 2011 international conference for high performance computing,*
1648 *networking, storage and analysis,* pages 1–32, 2011.
- 1649 29 D. El Baz, P. Spitéri, J.C. Miellou, and D. Gazen. Asynchronous iterative algorithms with
1650 flexible communication for non linear network flow problems. *Journal of Parallel and Distributed*
1651 *Computing,* 38:1–15, 1996.
- 1652 30 D. A. Beckingsale, J. Burmark, R. Hornung, H. Jones, W. Killian, A. J. Kunen, O. Pearce,
1653 P. Robinson, B. S. Ryujin, and T. R. Scogland. RAJA: Portable performance for large-scale
1654 scientific applications. In *2019 IEEE/ACM International Workshop on Performance, Portability*
1655 *and Productivity in HPC (P3HPC),* pages 71–81, Nov 2019. doi:10.1109/P3HPC49587.2019.
1656 00012.
- 1657 31 T. Benacchio, L. Bonaventura, M. Altenbernd, C.D. Cantwell, P.D. Düben, M. Gillard,
1658 L. Giraud, D. Göttsche, E. Raffin, K. Teranishi, and N. Wedi. Report on local data recovery
1659 approaches suitable for weather and climate prediction. FET-HPC ESCAPE-2 project
1660 deliverable, 2020. doi:10.2172/1607968.
- 1661 32 T. Benacchio, L. Bonaventura, M. Altenbernd, C.D. Cantwell, P.D. Düben, M. Gillard,
1662 L. Giraud, D. Göttsche, E. Raffin, K. Teranishi, and N. Wedi. Resilience and fault-tolerance
1663 in high-performance computing for numerical weather and climate prediction. *International*
1664 *Journal of High Performance Computing Applications,* 2020. Under revision.

- 1665 33 Anne Benoit, Aurélien Cavelan, Franck Cappello, Padma Raghavan, Yves Robert, and
1666 Hongyang Sun. Identifying the right replication level to detect and correct silent errors
1667 at scale. In *FTXS '17: Proceedings of the 2017 Workshop on Fault-Tolerance for HPC at*
1668 *Extreme Scale*, pages 31–38, 06 2017. doi:10.1145/3086157.3086162.
- 1669 34 Anne Benoit, Aurélien Cavelan, Yves Robert, and Hongyang Sun. Optimal resilience patterns
1670 to cope with fail-stop and silent errors. In *2016 IEEE International Parallel and Distributed*
1671 *Processing Symposium (IPDPS)*, pages 202–211. IEEE, 2016.
- 1672 35 Austin R. Benson, Sven Schmit, and Robert Schreiber. Silent error detection in numerical
1673 time-stepping schemes. *IJHPCA*, 29(4):403–421, 2015. arXiv:https://doi.org/10.1177/
1674 1094342014532297, doi:10.1177/1094342014532297.
- 1675 36 Marsha J. Berger, David L. George, Randall J. LeVeque, and Kyle T. Mandli. The GeoClaw
1676 software for depth-averaged flows with adaptive refinement. *Advances in Water Resources*,
1677 34(9):1195–1206, 2011.
- 1678 37 Eduardo Berrocal, Leonardo Bautista-Gomez, Sheng Di, Zhiling Lan, and Franck Cappello.
1679 Exploring partial replication to improve lightweight silent data corruption detection for HPC
1680 applications. In *Euro-Par 2016: Parallel Processing - 22nd International Conference on*
1681 *Parallel and Distributed Computing, Grenoble, France, August 24-26, 2016, Proceedings*, pages
1682 419–430, 2016. doi:10.1007/978-3-319-43659-3_31.
- 1683 38 Eduardo Berrocal, Leonardo Bautista-Gomez, Sheng Di, Zhiling Lan, and Franck Cappello.
1684 Toward general software level silent data corruption detection for parallel applications. *IEEE*
1685 *Trans. Parallel Distrib. Syst.*, 28(12):3642–3655, 2017. doi:10.1109/TPDS.2017.2735971.
- 1686 39 D. Bertsekas. Distributed asynchronous computation of fixed points. *Math. Programming*,
1687 27:107 – 120, 1983.
- 1688 40 D.P. Bertsekas and J.N. Tsitsiklis. *Parallel and Distributed Computation, Numerical Methods*.
1689 Prentice Hall, Englewood Cliffs N.J., 1989.
- 1690 41 Wesley Bland, Aurelien Bouteiller, Thomas Herault, Joshua Hursey, George Bosilca, and Jack J.
1691 Dongarra. An evaluation of user-level failure mitigation support in MPI. In *Proceedings of the*
1692 *19th European Conference on Recent Advances in the Message Passing Interface, EuroMPI'12*,
1693 pages 193–203, Berlin, Heidelberg, 2012. Springer-Verlag. doi:10.1007/978-3-642-33518-1_
1694 24.
- 1695 42 The OpenMP Architecture Review Boards. OpenMP (Open Multi-Processing), 2019. URL:
1696 <https://www.openmp.org/>.
- 1697 43 L. Bonaventura, F. Casella, L. Delpopolo Carciopolo, and A. Ranade. A self adjusting
1698 multirate algorithm for robust time discretization of partial differential equations. *Computers*
1699 *and Mathematics with Applications*, 79:2086–2098, 2020.
- 1700 44 G. Bosilca, R. Delmas, J. Dongarra, and J. Langou. Algorithm-based fault tolerance applied
1701 to high performance computing. *Journal of Parallel and Distributed Computing*, 69(4):410 –
1702 416, 2009.
- 1703 45 Patrick G. Bridges, Kurt B. Ferreira, Michael A. Heroux, and Mark Hoemmen. Fault-tolerant
1704 linear solvers via selective reliability, 2012. arXiv:1206.1390.
- 1705 46 S. Buchwald, Manuel Mohr, and Andreas Zwinkau. Malleable invasive applications. *CEUR*
1706 *Workshop Proceedings*, 1337:123–126, 01 2015.
- 1707 47 H.-J. Bungartz and S. Dirnstorfer. Multivariate quadrature on adaptive sparse grids. *Computing*,
1708 71(1):89–114, Aug 2003. doi:10.1007/s00607-003-0016-4.
- 1709 48 Hans-Joachim Bungartz and Michael Griebel. Sparse grids. *Acta Numerica*, 13:147–269, 2004.
1710 doi:10.1017/S0962492904000182.
- 1711 49 Jon Calhoun, Luke Olson, and Marc Snir. Flipit: An llvm based fault injector for hpc. In
1712 *Revised Selected Papers, Part I, of the Euro-Par 2014 International Workshops on Parallel*
1713 *Processing - Volume 8805*, page 547–558, Berlin, Heidelberg, 2014. Springer-Verlag. doi:
1714 10.1007/978-3-319-14325-5_47.

- 1715 50 Jon Calhoun, Luke Olson, Marc Snir, and William D. Gropp. Towards a more fault resilient
1716 multigrid solver. In *Proceedings of the Symposium on High Performance Computing*, San
1717 Diego, CA, USA, 2015. Society for Computer Simulation International.
- 1718 51 C. D. Cantwell and A. S. Nielsen. A minimally intrusive low-memory approach to resilience
1719 for existing transient solvers. *Journal of Scientific Computing*, 78(1):565–581, 2019.
- 1720 52 Franck Cappello, Al Geist, Bill Gropp, Laxmikant Kale, Bill Kramer, and Marc Snir. Toward
1721 exascale resilience. *The International Journal of High Performance Computing Applications*,
1722 23(4):374–388, 2009.
- 1723 53 L. Delpopolo Carciopolo, L. Bonaventura, A. Scotti, and L. Formaggia. A conservative implicit
1724 multirate method for hyperbolic problems. *Computational Geosciences*, 23:647–664, 2019.
- 1725 54 Ricolindo L. Carino, Ali Mohammed, and Florina M. Ciorba. Dynamic Loop Self-scheduling
1726 For Load Balancing (DLS4LB), 2020. URL: <https://github.com/unibas-dmi-hpc/DLS4LB>.
- 1727 55 Henri Casanova, Yves Robert, Frédéric Vivien, and Dounia Zaidouni. On the impact of process
1728 replication on executions of large-scale parallel applications with coordinated checkpointing.
1729 *Future Generation Comp. Syst.*, 51:7–19, 2015.
- 1730 56 Marc Casas, Wilfried N. Gansterer, and Elias Wimmer. Resilient gossip-inspired all-reduce al-
1731 gorithms for high-performance computing: Potential, limitations, and open questions. *IJHPCA*,
1732 33(2), 2019. doi:10.1177/1094342018762531.
- 1733 57 A. Cavelan, R. M. Cabezón, and Florina M. Ciorba. Detection of silent data corruptions
1734 in smoothed particle hydrodynamics simulations. In *Proceedings of the 19th IEEE/ACM*
1735 *International Symposium on Cluster, Cloud and Grid Computing (CCGrid 2019)*, May 2019.
- 1736 58 A. Cavelan and F. M. Ciorba. Algorithm-based fault tolerance for parallel stencil computa-
1737 tions. In *IEEE International Conference on Cluster Computing (Cluster 2019)*, Albuquerque,
1738 September 2019.
- 1739 59 M. Chau, T. Garcia, and P. Spiteri. Asynchronous Schwarz methods applied to constrained
1740 mechanical structures in grid environment. *Advances in Engineering Software*, 74:1 – 15, 2014.
- 1741 60 D. Chazan and W. Miranker. Chaotic relaxation. *Linear Algebra and its Applications*, 2(2):199 –
1742 222, 1969. URL: <http://www.sciencedirect.com/science/article/pii/0024379569900287>,
1743 doi:[https://doi.org/10.1016/0024-3795\(69\)90028-7](https://doi.org/10.1016/0024-3795(69)90028-7).
- 1744 61 S. Chen, G. Bronevetsky, M. Casas-Guix, and L. Peng. Comprehensive algorithmic resilience for
1745 numeric applications. Technical report, Lawrence Livermore National Lab.(LLNL), Livermore,
1746 CA (United States), 2013.
- 1747 62 Z. Chen and J. Dongarra. Algorithm-based fault tolerance for fail-stop failures. *IEEE*
1748 *Transactions on Parallel and Distributed Systems*, 19(12):1628–1641, 2008.
- 1749 63 Zizhong Chen. Online-ABFT: An online algorithm based fault tolerance scheme for soft error
1750 detection in iterative methods. *SIGPLAN Not.*, 48(8), February 2013. doi:10.1145/2517327.
1751 2442533.
- 1752 64 Andrew A. Chien, Pavan Balaji, Nan Dun, Aiman Fang, Hajime Fujita, Kamil Iskra, Zachary A.
1753 Rubenstein, Ziming Zheng, Jeff R. Hammond, Ignacio Laguna, D. Richards, Anshu Dubey,
1754 Brian van Straalen, Mark Hoemmen, Michael A. Heroux, Keita Teranishi, and Andrew R.
1755 Siegel. Exploring versioned distributed arrays for resilience in scientific applications. *IJHPCA*,
1756 31(6):564–590, 2017. doi:10.1177/1094342016664796.
- 1757 65 Edmond Chow, Hartwig Anzt, and Jack J. Dongarra. Asynchronous iterative algorithm for
1758 computing incomplete factorizations on GPUs. In Julian M. Kunkel and Thomas Ludwig,
1759 editors, *High Performance Computing - 30th International Conference, ISC High Performance*
1760 *2015, Frankfurt, Germany, July 12-16, 2015, Proceedings*, volume 9137 of *Lecture Notes in*
1761 *Computer Science*, pages 1–16. Springer, 2015. doi:10.1007/978-3-319-20119-1_1.
- 1762 66 Edmond Chow and Aftab Patel. Fine-grained parallel incomplete LU factorization. *SIAM J.*
1763 *Scientific Computing*, 37(2), 2015. doi:10.1137/140968896.
- 1764 67 Cluster File Systems, Inc. Lustre: A scalable, high-performance file system. White paper, Avail-
1765 able at [https://cse.buffalo.edu/faculty/tkosar/cse710/papers/lustre-whitepaper](https://cse.buffalo.edu/faculty/tkosar/cse710/papers/lustre-whitepaper.pdf).
1766 pdf, 2007.

- 1767 68 Evan Coleman, Masha Sosonkina, and Edmond Chow. Fault tolerant variants of the fine-
1768 grained parallel incomplete LU factorization. In Lukás Polok, Masha Sosonkina, William I.
1769 Thacker, and Josef Weinbub, editors, *Proceedings of the 25th High Performance Computing*
1770 *Symposium, Virginia Beach, VA, USA, April 23 - 26, 2017*, pages 15:1–15:12. ACM, 2017.
1771 URL: <http://dl.acm.org/citation.cfm?id=3108111>.
- 1772 69 Camille Coti, Thomas Herault, Pierre Lemarinier, Laurence Pilard, Ala Rezmerita, Eric
1773 Rodriguezb, and Franck Cappello. Blocking vs. non-blocking coordinated checkpointing for
1774 large-scale fault tolerant MPI. In *SC'06: Proceedings of the 2006 ACM/IEEE conference on*
1775 *Supercomputing*, pages 18–18. IEEE, 2006.
- 1776 70 S. P. Crago, D. I. Kang, M. Kang, R. Kost, K. Singh, J. Suh, and J. P. Walters. Programming
1777 models and development software for a space-based many-core processor. In *4th Int. Conf. on*
1778 *Space Mission Challenges for Information Technology*, pages 95–102. IEEE, 2011.
- 1779 71 C. Denis, P. de Oliveira Castro, and E. Petit. Verificarlo: checking floating point accuracy
1780 through Monte Carlo arithmetic. In *ARITH'23, Silicon Valley, USA*, Jul 2016.
- 1781 72 S. Di and F. Cappello. Fast error-bounded lossy HPC data compression with SZ. In *2016*
1782 *IEEE IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages
1783 730–739. IEEE, 2016.
- 1784 73 Sheng Di, Mohamed Slim Bouguerra, Leonardo Bautista-Gomez, and Franck Cappello. Opti-
1785 mization of multi-level checkpoint model for large scale HPC applications. In *2014 IEEE 28th*
1786 *International Parallel and Distributed Processing Symposium*, pages 1181–1190. IEEE, 2014.
- 1787 74 Sheng Di and Franck Cappello. Fast error-bounded lossy HPC data compression with SZ.
1788 In *2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages
1789 730–739. IEEE, 2016.
- 1790 75 Catello Di Martino, Zbigniew Kalbarczyk, Ravishankar K Iyer, Fabio Baccanico, Joseph Fullop,
1791 and William Kramer. Lessons learned from the analysis of system failures at petascale: The
1792 case of blue waters. In *2014 44th Annual IEEE/IFIP International Conference on Dependable*
1793 *Systems and Networks*, pages 610–621. IEEE, 2014. doi:10.1109/DSN.2014.62.
- 1794 76 Jose Monsalve Diaz, Swaroop Pophale, Kyle Friedline, Oscar Hernandez, David E Bernholdt,
1795 and Sunita Chandrasekaran. Evaluating support for openmp offload features. In *Proceedings*
1796 *of the 47th International Conference on Parallel Processing Companion*, pages 1–10, 2018.
- 1797 77 Edsger W Dijkstra. Self-stabilizing systems in spite of distributed control. *Communications*
1798 *of the ACM*, 17(11):643–644, 1974.
- 1799 78 Daniel Drzisga, Lorenz John, Ulrich Růde, Barbara Wohlmuth, and Walter Zulehner. On the
1800 analysis of block smoothers for saddle point problems. *SIAM Journal on Matrix Analysis*
1801 *and Applications*, 39(2):932–960, 2018. arXiv:<https://doi.org/10.1137/16M1106304>, doi:
1802 10.1137/16M1106304.
- 1803 79 P. Du, A. Bouteiller, G. Bosilca, T. Herault, and J. Dongarra. Algorithm-based fault tolerance
1804 for dense matrix factorizations. *ACM SIGPLAN notices*, 47(8):225–234, 2012.
- 1805 80 Alejandro Duran, Eduard Ayguadé, Rosa M. Badia, Jesús Labarta, Luis Martinell, Xavier
1806 Martorell, and Judit Planas. OmpSs: a proposal for programming heterogeneous multi-core
1807 architectures. *Parallel Processing Letters*, 21(2):173–193, 2011. URL: [http://dblp.uni-trier.](http://dblp.uni-trier.de/db/journals/ppl/pp121.html#DuranABLMMP11)
1808 [de/db/journals/ppl/pp121.html#DuranABLMMP11](http://dblp.uni-trier.de/db/journals/ppl/pp121.html#DuranABLMMP11).
- 1809 81 P. Eberhart, J. Brajard, P. Fortin, and F. Jézéquel. High performance numerical validation
1810 using stochastic arithmetic. *Reliable Computing*, 21:35–52, 2015.
- 1811 82 H. Carter Edwards, Christian R. Trott, and Daniel Sunderland. Kokkos: Enabling many-
1812 core performance portability through polymorphic memory access patterns. *Journal of*
1813 *Parallel and Distributed Computing*, 74(12):3202 – 3216, 2014. Domain-Specific Lan-
1814 guages and High-Level Frameworks for High-Performance Computing. URL: [http://www.](http://www.sciencedirect.com/science/article/pii/S0743731514001257)
1815 [sciencedirect.com/science/article/pii/S0743731514001257](http://www.sciencedirect.com/science/article/pii/S0743731514001257), doi:[https://doi.org/10.](https://doi.org/10.1016/j.jpdc.2014.07.003)
1816 [1016/j.jpdc.2014.07.003](https://doi.org/10.1016/j.jpdc.2014.07.003).

- 1817 83 Mireille El Haddad, José C. Garay, Frédéric Magoulès, and Daniel B. Szyld. Synchronous
1818 and asynchronous optimized Schwarz methods for one-way subdivision of bounded domains.
1819 *Numerical Linear Algebra and Applications*, 27:e2279, 2020. 30 pages.
- 1820 84 J. Elliott, M. Hoemmen, and F. Mueller. Evaluating the impact of SDC on the GMRES iterative
1821 solver. In *2014 IEEE 28th International Parallel and Distributed Processing Symposium*, pages
1822 1193–1202, May 2014. doi:10.1109/IPDPS.2014.123.
- 1823 85 J. Elliott, K. Kharbas, D. Fiala, F. Mueller, K. Ferreira, and C. Engelmann. Combining
1824 partial redundancy and checkpointing for HPC. In *2012 IEEE 32nd International Conference
1825 on Distributed Computing Systems (ICDCS)*, pages 615–626, June 2012. doi:10.1109/ICDCS.
1826 2012.56.
- 1827 86 James Elliott. *Resilient Iterative Linear Solvers Running Through Errors*. PhD thesis, North
1828 Carolina State University, 2015.
- 1829 87 Christian Engelmann and Thomas Naughton. Toward a performance/resilience tool for hard-
1830 ware/software co-design of high-performance computing systems. In *2013 42nd International
1831 Conference on Parallel Processing*, pages 960–969. IEEE, 2013.
- 1832 88 Christian Engelmann, Geoffroy R. Vallée, Thomas Naughton, and Stephen L. Scott. Proactive
1833 fault tolerance using preemptive migration. In *Proceedings of the 17th Euromicro International
1834 Conference on Parallel, Distributed, and network-based Processing (PDP) 2009*, pages 252–257,
1835 Weimar, Germany, February 18-20, 2009. IEEE Computer Society, Los Alamitos, CA, USA.
1836 doi:10.1109/PDP.2009.31.
- 1837 89 Christian Engelmann, Geoffroy R. Vallée, and Swaroop Pophale. Concepts for OpenMP target
1838 offload resilience. In *Lecture Notes in Computer Science: Proceedings of the 15th International
1839 Workshop on OpenMP (IWOMP) 2019*, volume 11718, pages 78–93, Auckland, New Zealand,
1840 September 11-13, 2019. Springer Verlag, Berlin, Germany. doi:10.1007/978-3-030-28596-8_
1841 6.
- 1842 90 Charbel Farhat and Francois-Xavier Roux. A method of finite element tearing and intercon-
1843 necting and its parallel solution algorithm. *International Journal for Numerical Methods in
1844 Engineering*, 32(6):1205–1227, 1991.
- 1845 91 David Fiala, Frank Mueller, Christian Engelmann, Kurt Ferreira, Ron Brightwell, and Rolf
1846 Riesen. Detection and correction of silent data corruption for large-scale high-performance com-
1847 puting. In *Proceedings of the 25th IEEE/ACM International Conference on High Performance
1848 Computing, Networking, Storage and Analysis (SC) 2012*, pages 78:1–78:12, Salt Lake City, UT,
1849 USA, November 10-16, 2012. ACM Press, New York, NY, USA. doi:10.1109/SC.2012.49.
- 1850 92 David Fiala, Frank Mueller, Kurt Ferreira, and Christian Engelmann. Mini-Ckpts: Surviving
1851 OS failures in persistent memory. In *Proceedings of the 30th ACM International Conference
1852 on Supercomputing (ICS) 2016*, pages 7:1–7:14, Istanbul, Turkey, June 1-3, 2016. ACM Press,
1853 New York, NY, USA. doi:10.1145/2925426.2926295.
- 1854 93 David Fiala, Frank Mueller, and Kurt B. Ferreira. Flipsphere: A software-based DRAM error
1855 detection and correction library for HPC. In *Proceedings of the 20th IEEE/ACM International
1856 Symposium on Distributed Simulation and Real Time Applications (DS-RT) 2016*, pages
1857 19–28, London, UK, September 21-23, 2016. IEEE Computer Society, Los Alamitos, CA, USA.
1858 doi:10.1109/DS-RT.2016.27.
- 1859 94 P.K. Fok. A linearly fourth order multirate Runge–Kutta method with error control. *Journal
1860 of Scientific Computing*, pages 1–19, 2015.
- 1861 95 MPI Forum. The Message Passing Interface standard, 2019. URL: [https://www.mpi-forum.
1862 org/](https://www.mpi-forum.org/).
- 1863 96 Michael Frechtling and Philip H. W. Leong. MCALIB: measuring sensitivity to rounding error
1864 with monte carlo programming. *ACM TOPLAS*, 37(2):1–25, 2015.
- 1865 97 Andreas Frommer and Daniel B. Szyld. On asynchronous iterations. *Journal of Computational
1866 and Applied Mathematics*, 123:201–216, 2000.
- 1867 98 François Févotte and Bruno Lathuilière. Debugging and optimization of HPC programs in
1868 mixed precision with the verrou tool. In *Computational Reproducibility at Exascale Workshop*

- 1869 (CRE2018), in conjunction with the International Conference on High Performance Computing,
1870 Networking, Storage and Analysis (SC18), Dallas, USA, 2018.
- 1871 **99** M. Gamell, D. S. Katz, H. Kolla, J. Chen, S. Klasky, and M. Parashar. Exploring automatic,
1872 online failure recovery for scientific applications at extreme scales. In *SC '14: Proceedings*
1873 *of the International Conference for High Performance Computing, Networking, Storage and*
1874 *Analysis*, pages 895–906, Nov 2014. doi:10.1109/SC.2014.78.
- 1875 **100** Marc Gamell, Daniel S. Katz, Keita Teranishi, Michael A. Heroux, Rob F. Van der Wijngaart,
1876 Timothy G. Mattson, and Manish Parashar. Evaluating online global recovery with Fenix
1877 using application-aware in-memory checkpointing techniques. In *ICPP Workshops*, pages
1878 346–355. IEEE Computer Society, 2016. URL: [http://dblp.uni-trier.de/db/conf/icppw/
1879 icppw2016.html#GamellKTHWMP16](http://dblp.uni-trier.de/db/conf/icppw/icppw2016.html#GamellKTHWMP16).
- 1880 **101** Wilfried N. Gansterer, Gerhard Niederbrucker, Hana Straková, and Stefan Schulze Grotthoff.
1881 Robust distributed orthogonalization based on randomized aggregation. In Vassil N. Alexan-
1882 drov, Al Geist, and Jack J. Dongarra, editors, *Proceedings of the second workshop on Scalable*
1883 *algorithms for large-scale systems, ScalA@SC 2011, Seattle, WA, USA, November 14, 2011*,
1884 pages 7–10. ACM, 2011. doi:10.1145/2133173.2133177.
- 1885 **102** Wilfried N. Gansterer, Gerhard Niederbrucker, Hana Straková, and Stefan Schulze Grotthoff.
1886 Scalable and fault tolerant orthogonalization based on randomized distributed data aggregation.
1887 *J. Comput. Sci.*, 4(6):480–488, 2013. doi:10.1016/j.jocs.2013.01.006.
- 1888 **103** José C. Garay, Frédéric Magoulès, and Daniel B. Szyld. Synchronous and asynchronous
1889 optimized Schwarz method for Poisson’s equation in rectangular domains. Technical Report
1890 17-10-18, Department of Mathematics, Temple University, October 2017. Revised April 2018.
- 1891 **104** Jochen Garcke. Regression with the optimised combination technique. In *Proceedings of the*
1892 *23rd international conference on Machine learning*, pages 321–328. ACM Press, 2006.
- 1893 **105** Jochen Garcke. A dimension adaptive sparse grid combination technique for machine learning.
1894 *ANZIAM Journal*, 48:725–740, 2007.
- 1895 **106** C.W. Gear and D.R. Wells. Multirate linear multistep methods. *BIT*, 24:484–502, 1984.
- 1896 **107** Giorgis Georgakoudis, Ignacio Laguna, Dimitrios S. Nikolopoulos, and Martin Schulz. Refine:
1897 Realistic fault injection via compiler-based instrumentation for accuracy, portability and speed.
1898 In *Proceedings of the International Conference for High Performance Computing, Networking,*
1899 *Storage and Analysis, SC '17, New York, NY, USA, 2017*. Association for Computing Machinery.
1900 doi:10.1145/3126908.3126972.
- 1901 **108** Cijo George and Sathish S. Vadhiyar. ADFT: An adaptive framework for fault tolerance on
1902 large scale systems using application malleability. *Procedia Computer Science*, 9:166 – 175,
1903 2012.
- 1904 **109** T. Gerstner and M. Griebel. Dimension-adaptive tensor-product quadrature. *Computing*,
1905 71(1):65–87, Aug 2003. doi:10.1007/s00607-003-0015-5.
- 1906 **110** Thomas Gerstner and Michael Griebel. Numerical integration using sparse grids. *Numerical*
1907 *Algorithms*, 18(3):209, Jan 1998. doi:10.1023/A:1019129717644.
- 1908 **111** Roberto Gioiosa, Jose Carlos Sancho, Song Jiang, and Fabrizio Petrini. Transparent, incremen-
1909 tal checkpointing at kernel level: a foundation for fault tolerance for parallel computers. In
1910 *SC'05: Proceedings of the 2005 ACM/IEEE conference on Supercomputing*, pages 9–9. IEEE,
1911 2005.
- 1912 **112** Christian Glusa, Erik G. Boman, Edmond Chow, Sivasankaran Rajamanickam, and Daniel B.
1913 Szyld. Scalable asynchronous domain decomposition solvers. Technical Report 19-10-11,
1914 Department of Mathematics, Temple University, October 2019.
- 1915 **113** Christian Glusa, Paritosh Ramanan, Erik G. Boman, Edmond Chow, and Sivasankaran
1916 Rajamanickam. Asynchronous one-level and two-level domain decomposition solvers. *CoRR*,
1917 abs/1808.08172, 2018. URL: <http://arxiv.org/abs/1808.08172>, arXiv:1808.08172.
- 1918 **114** Björn Gmeiner, Ulrich Rüde, Holger Stengel, Christian Waluga, and Barbara Wohlmuth.
1919 Towards textbook efficiency for parallel multigrid. *Numerical Mathematics: Theory, Methods*
1920 *and Applications*, 8(1):22–46, 2015.

- 1921 115 Dominik Göttsche, Mirco Altenbernd, and Dirk Ribbrock. Fault-tolerant finite-element
1922 multigrid algorithms with hierarchically compressed asynchronous checkpointing. *Parallel*
1923 *Computing*, 49:117–135, October 2015. doi:10.1016/j.parco.2015.07.003.
- 1924 116 Michael Griebel and Peter Oswald. Stochastic subspace correction methods and fault tolerance.
1925 *Mathematics of Computation*, 89(321):279–312, 2020.
- 1926 117 Michael Griebel, Michael Schneider, and Christoph Zenger. A combination technique for the
1927 solution of sparse grid problems. In *Iterative Methods in Lin. Alg.*, pages 263–281, 1992.
- 1928 118 The Khronos Group. OpenCL (Open Computing Language), 2020. URL: <https://www.khronos.org/opencv/>.
- 1930 119 Ray Grout, Hemanth Kolla, Michael Minion, and John Bell. Achieving algorithmic resilience
1931 for temporal integration through spectral deferred corrections. *Communications in Applied*
1932 *Mathematics and Computational Science*, 12(1):25–50, 2017.
- 1933 120 Q. Guan, N. Debardeleben, S. Blanchard, and S. Fu. F-SEFI: A fine-grained soft error fault
1934 injection tool for profiling application vulnerability. In *2014 IEEE 28th International Parallel*
1935 *and Distributed Processing Symposium*, pages 1245–1254, May 2014. doi:10.1109/IPDPS.
1936 2014.128.
- 1937 121 Pierre-Louis Guhur, Hong Zhang, Tom Peterka, Emil Constantinescu, and Franck Cappello.
1938 Lightweight and accurate silent data corruption detection in ordinary differential equation
1939 solvers. In *European Conference on Parallel Processing*, pages 644–656. Springer, 2016.
- 1940 122 E. Hairer, S.P. Nørsett, and G. Wanner. *Solving Ordinary Differential Equations I: Nonstiff*
1941 *Problems*. Springer-Verlag, Berlin Heidelberg, 3rd corr. edition, 2008.
- 1942 123 Brendan Harding et al. Fault tolerant computation with the sparse grid combination technique.
1943 *SIAM Journal on Scient. Comp.*, 37(3):C331–C353, 2015.
- 1944 124 Brendan Harding and Markus Hegland. Robust solutions to PDEs with multiple grids. In
1945 Jochen Garcke and Dirk Pflüger, editors, *Sparse Grids and Applications - Munich 2012 SE*,
1946 volume 97 of *Lecture Notes in Computational Science and Engineering*, pages 171–193. Springer
1947 International Publishing, 2014.
- 1948 125 Paul H. Hargrove and Jason C. Duell. Berkeley Lab Checkpoint/Restart (BLCR) for Linux
1949 clusters. In *Journal of Physics: Proceedings of the Scientific Discovery through Advanced*
1950 *Computing Program (SciDAC) Conference 2006*, volume 46, pages 494–499, Denver, CO, USA,
1951 June 25–29, 2006. Institute of Physics Publishing, Bristol, UK. doi:10.1088/1742-6596/46/
1952 1/067.
- 1953 126 Siva Kumar Sastry Hari, Timothy Tsai, Mark Stephenson, Stephen W. Keckler, and Joel S.
1954 Emer. SASSIFI: an architecture-level fault injection tool for GPU application resilience
1955 evaluation. In *2017 IEEE International Symposium on Performance Analysis of Systems*
1956 *and Software, ISPASS 2017, Santa Rosa, CA, USA, April 24–25, 2017*, pages 249–258. IEEE
1957 Computer Society, 2017. doi:10.1109/ISPASS.2017.7975296.
- 1958 127 Amin Hassani, Anthony Skjellum, and Ron Brightwell. Design and evaluation of FA-MPI, a
1959 transactional resilience scheme for non-blocking MPI. In *44th Annual IEEE/IFIP International*
1960 *Conference on Dependable Systems and Networks, DSN 2014, Atlanta, GA, USA, June 23–26,*
1961 *2014*, pages 750–755. IEEE Computer Society, 2014. doi:10.1109/DSN.2014.78.
- 1962 128 Xubin (Ben) He, Li Ou, Christian Engelmann, Xin Chen, and Stephen L. Scott. Symmetric
1963 active/active metadata service for high availability parallel file systems. *Journal of Parallel*
1964 *and Distributed Computing (JPDC)*, 69(12):961–973, December 2009. doi:10.1016/j.jpdc.
1965 2009.08.004.
- 1966 129 Mario Heene, Alfredo Parra Hinojosa, Michael Obersteiner, Hans-Joachim Bungartz, and
1967 Dirk Pflüger. EXAHD: An exa-scalable two-level sparse grid approach for higher-dimensional
1968 problems in plasma physics and beyond. In *High Performance Computing in Science and*
1969 *Engineering'17*, pages 513–529. Springer, 2018.
- 1970 130 Mario Heene, Alfredo Parra Hinojosa, Hans-Joachim Bungartz, and Dirk Pflüger. A massively-
1971 parallel, fault-tolerant solver for high-dimensional PDEs. In *Euro-Par 2016: Parallel Processing*
1972 *Workshops*, 2016.

- 1973 131 Thomas Herault and Yves Robert, editors. *Fault-Tolerance Techniques for High-Performance*
1974 *Computing*. Springer Verlag, 2015.
- 1975 132 Magnus R. Hestenes and Eduard Stiefel. Methods of conjugate gradients for solving linear
1976 systems. *Journal of Research of the National Bureau of Standards*, 46(6):409–436, December
1977 1952.
- 1978 133 N.J. Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and
1979 Applied Mathematics, Philadelphia, PA, USA, 1st edition, 1996.
- 1980 134 Alfredo Parra Hinojosa, Brendan Harding, Markus Hegland, and Hans-Joachim Bungartz.
1981 Handling silent data corruption with the sparse grid combination technique. In *Software for*
1982 *Exascale Computing-SPPEXA 2013-2015*, pages 187–208. Springer, 2016.
- 1983 135 Mark Hoemmen and Michael Allen Heroux. Fault-tolerant iterative methods via selec-
1984 tive reliability. *Proceedings of the IEEE/ACM conference on supercomputing (SC'11)*,
1985 2011. URL: [http://www.researchgate.net/publication/228940928_Fault-Tolerant_](http://www.researchgate.net/publication/228940928_Fault-Tolerant_Iterative_Methods_via_Selective_Reliability/file/79e41508060305e4ad.pdf)
1986 [Iterative_Methods_via_Selective_Reliability/file/79e41508060305e4ad.pdf](http://www.researchgate.net/publication/228940928_Fault-Tolerant_Iterative_Methods_via_Selective_Reliability/file/79e41508060305e4ad.pdf).
- 1987 136 K.-H. Huang and J.A. Abraham. Algorithm-based fault tolerance for matrix operations. *IEEE*
1988 *Transactions on Computers*, 100(6):518–528, 1984.
- 1989 137 Kuang-hua Huang and Jacob a Abraham. Algorithm-based fault tolerance for matrix operations.
1990 *IEEE Transactions on Computers*, c(6):518–528, 1984. doi:10.1109/TC.1984.1676475.
- 1991 138 Markus Huber, Björn Gmeiner, Ulrich Rüde, and Barbara Wohlmuth. Resilience for massively
1992 parallel multigrid solvers. *SIAM Journal on Scientific Computing*, 38(5):S217–S239, 2016.
- 1993 139 Markus Huber, Ulrich Rüde, and Barbara Wohlmuth. Adaptive control in roll-forward recovery
1994 for extreme scale multigrid. *The International Journal of High Performance Computing*
1995 *Applications*, 33(5):817–837, 2019.
- 1996 140 S. Hukerikar, K. Teranishi, P. C. Diniz, and R. F. Lucas. An evaluation of lazy fault detection
1997 based on adaptive redundant multithreading. In *2014 IEEE High Performance Extreme*
1998 *Computing Conference (HPEC)*, pages 1–6, 2014.
- 1999 141 Saurabh Hukerikar and Christian Engelmann. A pattern language for high-performance
2000 computing resilience. In *Proceedings of the 22nd European Conference on Pattern Languages*
2001 *of Programs (EuroPLoP) 2017*, pages 12:1–12:16, Kloster Irsee, Germany, July 12-16, 2017.
2002 ACM Press, New York, NY, USA. doi:10.1145/3147704.3147718.
- 2003 142 Saurabh Hukerikar and Christian Engelmann. Resilience design patterns: A structured
2004 approach to resilience at extreme scale (version 1.2). Technical Report ORNL/TM-2017/745,
2005 Oak Ridge National Laboratory, Oak Ridge, TN, USA, August 2017. doi:10.2172/1436045.
- 2006 143 Frank Hülsemann, Markus Kowarschik, Marcus Mohr, and Ulrich Rüde. Parallel geometric
2007 multigrid. In *Numerical Solution of Partial Differential Equations on Parallel Computers*,
2008 pages 165–208. Springer, 2006.
- 2009 144 R. K. Iyer, N. M. Nakka, Z. T. Kalbarczyk, and S. Mitra. Recent advances and new
2010 avenues in hardware-level reliability support. *IEEE Micro*, 25(6):18–29, Nov 2005. doi:
2011 10.1109/MM.2005.119.
- 2012 145 John D. Jakeman and Stephen G. Roberts. Local and dimension adaptive sparse grid
2013 interpolation and quadrature. *arXiv preprint arXiv:1110.0010*, September, 2011. arXiv:
2014 1110.0010v1.
- 2015 146 L. Jaulmes, M. Casas, M. Moretó ans E. Ayguadé, J. Labarta, and M. Valero. Exploiting
2016 asynchrony from exact forward recovery for detected and uncorrected errors in iterative solvers.
2017 In *SC '15: Proceedings of the International Conference for High Performance Computing,*
2018 *Networking, Storage and Analysis*, Austin, TX, 2015.
- 2019 147 Y. Jia, G. Bosilca, P. Luszczek, and J. J. Dongarra. Parallel reduction to Hessenberg form
2020 with algorithm-based fault tolerance. In *Proceedings of the International Conference on High*
2021 *Performance Computing, Networking, Storage and Analysis*, SC '13, New York, NY, USA,
2022 2013. Association for Computing Machinery. doi:10.1145/2503210.2503249.

- 2023 148 A. Jin, J. Jiang, J. Hu, and J. Lou. A pin-based dynamic software fault injection system. In
2024 2008 *The 9th International Conference for Young Computer Scientists*, pages 2160–2167, Nov
2025 2008. doi:10.1109/ICYCS.2008.329.
- 2026 149 U. Kabir and D. Goswami. An ABFT scheme based on communication characteristics. In
2027 2016 *IEEE International Conference on Cluster Computing (CLUSTER)*, pages 515–523, Sep.
2028 2016. doi:10.1109/CLUSTER.2016.68.
- 2029 150 Hartmut Kaiser, Thomas Heller, Bryce Adelstein-Lelbach, Adrian Serio, and Dietmar Fey.
2030 HPX: A task based programming model in a global address space. In *Proceedings of the 8th*
2031 *International Conference on Partitioned Global Address Space Programming Models, PGAS'14*,
2032 New York, NY, USA, 2014. Association for Computing Machinery. doi:10.1145/2676870.
2033 2676883.
- 2034 151 Laxmikant V. Kale and Sanjeev Krishnan. Charm++: A portable concurrent object oriented
2035 system based on C++. In *Proceedings of the Eighth Annual Conference on Object-Oriented*
2036 *Programming Systems, Languages, and Applications, OOPSLA'93*, page 91–108, New York,
2037 NY, USA, 1993. Association for Computing Machinery. doi:10.1145/165854.165874.
- 2038 152 George Karniadakis and Spencer Sherwin. *Spectral/hp element methods for computational*
2039 *fluid dynamics*. Oxford University Press, 2013.
- 2040 153 Kai Keller and Leonardo Bautista-Gomez. Application-level differential checkpointing for HPC
2041 applications with dynamic datasets. In *19th IEEE/ACM International Symposium on Cluster,*
2042 *Cloud and Grid Computing, CCGRID 2019, Larnaca, Cyprus, May 14-17, 2019*, pages 52–61.
2043 IEEE, 2019. doi:10.1109/CCGRID.2019.00015.
- 2044 154 Jungwon Kim, Seyong Lee, and Jeffrey S. Vetter. PapyrusKV: A high-performance parallel key-
2045 value store for distributed NVM architectures. In *Proceedings of the International Conference*
2046 *for High Performance Computing, Networking, Storage and Analysis*, New York, NY, USA,
2047 2017. Association for Computing Machinery. doi:10.1145/3126908.3126943.
- 2048 155 Axel Klawonn, Martin J. Kühn, and Oliver Rheinbach. Parallel adaptive FETI-DP using
2049 lightweight asynchronous dynamic load balancing. *International Journal for Numerical*
2050 *Methods in Engineering*, 121(4):621–643, 2020. URL: [https://onlinelibrary.wiley.com/](https://onlinelibrary.wiley.com/doi/abs/10.1002/nme.6237)
2051 [doi/pdf/10.1002/](https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.6237)
2052 [nme.6237](https://onlinelibrary.wiley.com/doi/pdf/10.1002/nme.6237), doi:10.1002/nme.6237.
- 2053 156 Nils Kohl, Johannes Hötzer, Florian Schornbaum, Martin Bauer, Christian Godenschwager,
2054 Harald Köstler, Britta Nestler, and Ulrich Rüde. A scalable and extensible checkpointing
2055 scheme for massively parallel simulations. *The International Journal of High Performance*
2056 *Computing Applications*, 33(4):571–589, 2019.
- 2057 157 Giorgos Kollias, Efstratios Gallopoulos, and Daniel B. Szyld. Asynchronous iterative com-
2058 putations with Web information retrieval structures: The PageRank case. In G.R. Joubert,
2059 W.E. Nagel, F.J. Peters, O. Plata, P. Tirado, and E. Zapata, editors, *Parallel Computing:*
2060 *Current and Future Issues of High-End Computing (Proceedings of the International Confer-*
2061 *ence Parco05)*, volume 33 of *NIC Series*, pages 309–316, Jülich, Germany, 2006. John von
2062 Neumann-Institut für Computing (NIC).
- 2063 158 U.W. Kulisch. *Advanced Arithmetic for the Digital Computer*. Springer, Wien, 2002.
- 2064 159 Ignacio Laguna, David Richards, Todd Gamblin, Martin Schulz, Bronis Supinski, Kathryn
2065 Mohror, and Howard Pritchard. Evaluating and extending user-level fault tolerance in MPI
2066 applications. *The International Journal of High Performance Computing Applications*, 30, 01
2067 2016. doi:10.1177/1094342015623623.
- 2068 160 J.D. Lambert. *Numerical Methods for Ordinary Differential Systems*. Wiley, Chirchester,
2069 England, 1991.
- 2070 161 Julien Langou, Zizhong Chen, George Bosilca, and Jack Dongarra. Recovery patterns for
2071 iterative methods in a parallel unstable environment. *SIAM J. Sci. Comput.*, 30:102–116,
2072 November 2007. doi:10.1137/040620394.
- 2073 162 Christian Lengauer, Sven Apel, Mathias Bolten, Shigeru Chiba, Ulrich Rüde, Jürgen Teich,
2074 Armin Größlinger, Frank Hannig, Harald Köstler, Lisa Claus, Alexander Grebhahn, Stefan

- 2075 Groth, Stefan Kronawitter, Sebastian Kuckuk, Hannah Rittich, Christian Schmitt, and Jonas
2076 Schmitt. ExaStencils – Advanced Multigrid Solver Generation. In Hans-Joachim Bungartz,
2077 Severin Reiz, Philipp Neumann, Benjamin Uekermann, and Wolfgang Nagel, editors, *Software
2078 for Exascale Computing – SPPEXA 2016-2019*, Lecture Notes in Computer Science and Engi-
2079 neering. Springer, 2020. URL: [https://www12.cs.fau.de/downloads/hannig/publications/
2080 ExaStencils_Advanced_Multigrid_Solver_Generation.pdf](https://www12.cs.fau.de/downloads/hannig/publications/ExaStencils_Advanced_Multigrid_Solver_Generation.pdf).
- 2081 **163** Randall J. LeVeque, David L. George, and Marsha J. Berger. Tsunami modelling with
2082 adaptively refined finite volume methods. *Acta Numerica*, 20:211–289, 2011.
- 2083 **164** Markus Levonyak, Christina Pacher, and Wilfried N. Gansterer. Scalable resilience against
2084 node failures for communication-hiding preconditioned conjugate gradient and conjugate resid-
2085 ual methods. In *Proceedings of the 2020 SIAM Conference on Parallel Processing for Scientific
2086 Computing*, pages 81–92. SIAM, 2020. URL: [https://epubs.siam.org/doi/abs/10.1137/1.
2087 9781611976137.8](https://epubs.siam.org/doi/abs/10.1137/1.9781611976137.8), arXiv:<https://epubs.siam.org/doi/pdf/10.1137/1.9781611976137.8>,
2088 doi:10.1137/1.9781611976137.8.
- 2089 **165** Scott Levy, Bryan Topp, Kurt B Ferreira, Dorian Arnold, Torsten Hoefler, and Patrick Widener.
2090 Using simulation to evaluate the performance of resilience strategies at scale. In *International
2091 Workshop on Performance Modeling, Benchmarking and Simulation of High Performance
2092 Computer Systems*, pages 91–114. Springer, 2013.
- 2093 **166** X. Liang, S. Di, D. Tao, Si. Li, Sh. Li, H. Guo, Z. Chen, and F. Cappello. Error-controlled
2094 lossy compression optimized for high compression ratios of scientific datasets. *2018 IEEE
2095 International Conference on Big Data (Big Data)*, pages 438–447, 2018.
- 2096 **167** Xin Liang, Jieyang Chen, Dingwen Tao, Sihuan Li, Panruo Wu, Hongbo Li, Kaiming Ouyang,
2097 Yuanlai Liu, Fengguang Song, and Zizhong Chen. Correcting soft errors online in fast Fourier
2098 transform. In *Proceedings of the International Conference for High Performance Computing,
2099 Networking, Storage and Analysis*, New York, NY, USA, 2017. Association for Computing
2100 Machinery. doi:10.1145/3126908.3126915.
- 2101 **168** J. Lidman, D. J. Quinlan, C. Liao, and S. A. McKee. ROSE::FTTransform - a source-to-source
2102 translation framework for exascale fault-tolerance research. In *IEEE/IFIP International
2103 Conference on Dependable Systems and Networks Workshops (DSN 2012)*, pages 1–6, June
2104 2012. doi:10.1109/DSNW.2012.6264672.
- 2105 **169** J. Liesen and Z. Strakoš. *Krylov Subspace Methods*. Numerical Mathematics and Scientific
2106 Computation. Oxford University Press, 2013.
- 2107 **170** S. Lin and P. Chen. A simd-based software fault tolerance for arm processors. In *2017
2108 International Conference on Applied System Innovation (ICASI)*, pages 910–913, 2017.
- 2109 **171** Romain Lion. Tolérance aux pannes dans l’exécution distribuée de graphes de tâches. In
2110 *Conférence d’informatique en Parallélisme, Architecture et Système*, Anglet, France, June 2019.
2111 URL: <https://hal.inria.fr/hal-02296118>.
- 2112 **172** Romain Lion and Samuel Thibault. From tasks graphs to asynchronous distributed check-
2113 pointing with local restart. In *Fault Tolerance for HPC at eXtreme Scale (FTXS) Workshop,
2114 2020*.
- 2115 **173** N. Losada, G. Bosilca, A. Bouteiller, P. González, and M. J. Martín. Local rollback for resilient
2116 MPI applications with application-level checkpointing and message logging. *Future Generation
2117 Computer Systems*, 91:450–464, 2019.
- 2118 **174** Kaoutar Maghraoui, Travis Desell, Boleslaw Szymanski, and Carlos Varela. Dynamic malleabil-
2119 ity in iterative MPI applications. In *In: 7th International Symposium on Cluster Computing
2120 and the Grid*, pages 591–598, 05 2007. doi:10.1109/CCGRID.2007.45.
- 2121 **175** Frédéric Magoulès, Daniel B. Szyld, and Cédric Venet. Asynchronous optimized Schwarz
2122 methods with and without overlap. *Numerische Mathematik*, 137:199–227, 2017.
- 2123 **176** Tim Mattson, Romain Cledat, Vincent Cave, Vivek Sarkar, Zoran Budimlic, Sanjay Chatterjee,
2124 Josh Fryman, Ivan Ganey, Robin Knauerhase, Min Lee, Benoit Meister, Brian Nickerson,
2125 Nick Pepperling, Bala Seshasayee, Sagnak Tasirlar, Justin Teller, and Nick Vrvilo. The open
2126 community runtime: A runtime system for extreme scale computing. In *2016 IEEE High*

- 2127 *Performance Extreme Computing Conference (HPEC)*, pages 1–7, 09 2016. doi:10.1109/
2128 HPEC.2016.7761580.
- 2129 **177** Harshitha Menon and Kathryn Mohror. DisCVar: Discovering critical variables using
2130 algorithmic differentiation for transient faults. *SIGPLAN Not.*, 53(1), February 2018.
2131 doi:10.1145/3200691.3178502.
- 2132 **178** G. Meurant and Z. Strakoš. The Lanczos and Conjugate Gradient algorithms in finite precision
2133 arithmetic. *Acta Numerica*, 15:471–542, 2006.
- 2134 **179** Amitabh Mishra and Prithviraj Banerjee. An algorithm-based error detection scheme for
2135 the multigrid method. *IEEE Trans. Comput.*, 52(9), September 2003. doi:10.1109/TC.2003.
2136 1228507.
- 2137 **180** A. Mohammed, A. Cavelan, and F. M. Ciorba. rDLB: A novel approach for robust dynamic
2138 load balancing of scientific applications with independent tasks. In *Proceedings of the 2019*
2139 *International Conference on High Performance Computing & Simulation (HPCS 2019)*, Dublin,
2140 July 2019.
- 2141 **181** Michael Moldaschl, Karl E. Prikopa, and Wilfried N. Gansterer. Fault tolerant communication-
2142 optimal 2.5D matrix multiplication. *J. Parallel Distributed Comput.*, 104:179–190, 2017.
2143 doi:10.1016/j.jpdc.2017.01.022.
- 2144 **182** A. Moody, G. Bronevetsky, K. Mohror, and B. R. d. Supinski. Design, modeling, and
2145 evaluation of a scalable multi-level checkpointing system. In *SC '10: Proceedings of the 2010*
2146 *ACM/IEEE International Conference for High Performance Computing, Networking, Storage*
2147 *and Analysis*, pages 1–11, Nov 2010. doi:10.1109/SC.2010.18.
- 2148 **183** Adam Moody, Greg Bronevetsky, Kathryn Mohror, and Bronis R. de Supinski. Design,
2149 modeling, and evaluation of a scalable multi-level checkpointing system. In *Proceedings of the*
2150 *2010 ACM/IEEE International Conference for High Performance Computing, Networking,*
2151 *Storage and Analysis*, SC '10, page 1–11, USA, 2010. IEEE Computer Society. doi:10.1109/
2152 SC.2010.18.
- 2153 **184** K. Morris, F. Rizzi, B. Cook, P. Mycek, O. LeMaitre, O. M. Knio, K. Sargsyan, K. Dahlgren,
2154 and B. J. Debusschere. Performance scaling variability and energy analysis for a resilient
2155 ULFM-based PDE solver. In *2016 7th Workshop on Latest Advances in Scalable Algorithms*
2156 *for Large-Scale Systems (ScalA)*, pages 41–48, Nov 2016. doi:10.1109/ScalA.2016.010.
- 2157 **185** Andreas Müller, Jörn Behrens, Francis X Giraldo, and Volkmar Wirth. Comparison between
2158 adaptive and uniform discontinuous Galerkin simulations in dry 2D bubble experiments.
2159 *Journal of Computational Physics*, 235:371–393, 2013.
- 2160 **186** Paul Mycek, Francesco Rizzi, Olivier Le Maître, Khachik Sargsyan, Karla Morris, Cosmin Safta,
2161 Bert Debusschere, and Omar Knio. Discrete a priori bounds for the detection of corrupted PDE
2162 solutions in exascale computations. *SIAM Journal on Scientific Computing*, 39(1):C1–C28,
2163 2017. arXiv:<https://doi.org/10.1137/15M1051786>, doi:10.1137/15M1051786.
- 2164 **187** Arun B. Nagarajan, Frank Mueller, Christian Engelmann, and Stephen L. Scott. Proactive
2165 fault tolerance for HPC with Xen virtualization. In *Proceedings of the 21st ACM International*
2166 *Conference on Supercomputing (ICS) 2007*, pages 23–32, Seattle, WA, USA, June 16-20, 2007.
2167 ACM Press, New York, NY, USA. doi:10.1145/1274971.1274978.
- 2168 **188** X. Ni and L. V. Kale. FlipBack: Automatic targeted protection against silent data corruption.
2169 In *Conference for High Performance Computing, Networking, Storage and Analysis (SC)*,
2170 pages 335–346, Nov 2016. doi:10.1109/SC.2016.28.
- 2171 **189** Bogdan Nicolae, Adam Moody, Elsa Gonsiorowski, Kathryn Mohror, and Franck Cappello.
2172 VeloC: Towards high performance adaptive asynchronous checkpointing at large scale. In
2173 *IPDPS'19: The 2019 IEEE International Parallel and Distributed Processing Symposium*,
2174 pages 911–920, Rio de Janeiro, Brazil, May 2019. URL: [https://hal.archives-ouvertes.](https://hal.archives-ouvertes.fr/hal-02184203)
2175 [fr/hal-02184203](https://hal.archives-ouvertes.fr/hal-02184203).
- 2176 **190** Gerhard Niederbrucker and Wilfried N. Gansterer. Robust gossip-based aggregation: A
2177 practical point of view. In Peter Sanders and Norbert Zeh, editors, *Proceedings of the 15th*

- 2178 *Meeting on Algorithm Engineering and Experiments, ALLENEX 2013, New Orleans, Louisiana,*
2179 *USA, January 7, 2013*, pages 133–147. SIAM, 2013. doi:10.1137/1.9781611972931.12.
- 2180 **191** Allan S. Nielsen and Jan S. Hesthaven. Fault tolerance in the Parareal method. In *Proceedings*
2181 *of the ACM Workshop on Fault-Tolerance for HPC at Extreme Scale, FTXS '16*, pages 1–8,
2182 New York, NY, USA, 2016. ACM. URL: <http://dx.doi.org/10.1145/2909428.2909431>,
2183 doi:10.1145/2909428.2909431.
- 2184 **192** Michael Obersteiner, Alfredo Parra Hinojosa, Mario Heene, Hans-Joachim Bungartz, and
2185 Dirk Pflüger. A highly scalable, algorithm-based fault-tolerant solver for gyrokinetic plasma
2186 simulations. In *Proceedings of the 8th Workshop on Latest Advances in Scalable Algorithms*
2187 *for Large-Scale Systems*, pages 1–8, 2017.
- 2188 **193** N. Oh, P. P. Shirvani, and E. J. McCluskey. Error detection by duplicated instructions in
2189 super-scalar processors. *IEEE Transactions on Reliability*, 51(1):63–75, March 2002. doi:
2190 10.1109/24.994913.
- 2191 **194** C. Pachajoa, C. Pacher, and W. N. Gansterer. Node-failure-resistant preconditioned conjugate
2192 gradient method without replacement nodes. In *2019 IEEE/ACM 9th Workshop on Fault*
2193 *Tolerance for HPC at eXtreme Scale (FTXS)*, pages 31–40, Nov 2019. doi:10.1109/FTXS49593.
2194 2019.00009.
- 2195 **195** Carlos Pachajoa and Wilfried N. Gansterer. On the resilience of conjugate gradient and
2196 multigrid methods to node failures. In Dora Blanco Heras, Luc Bougé, Gabriele Mencagli,
2197 Emmanuel Jeannot, Rizos Sakellariou, Rosa M. Badia, Jorge G. Barbosa, Laura Ricci,
2198 Stephen L. Scott, Stefan Lankes, and Josef Weidendorfer, editors, *Euro-Par 2017: Parallel*
2199 *Processing Workshops - Euro-Par 2017 International Workshops, Santiago de Compostela,*
2200 *Spain, August 28-29, 2017, Revised Selected Papers*, volume 10659 of *Lecture Notes in Computer*
2201 *Science*, pages 569–580. Springer, 2017. doi:10.1007/978-3-319-75178-8_46.
- 2202 **196** Carlos Pachajoa, Markus Levonyak, and Wilfried N. Gansterer. Extending and evaluating
2203 fault-tolerant preconditioned conjugate gradient methods. In *IEEE/ACM 8th Workshop on*
2204 *Fault Tolerance for HPC at eXtreme Scale, FTXS@SC 2018, Dallas, TX, USA, November 16,*
2205 *2018*, pages 49–58. IEEE, 2018. doi:10.1109/FTXS.2018.00009.
- 2206 **197** Carlos Pachajoa, Markus Levonyak, Wilfried N. Gansterer, and Jesper Larsson Träff. How to
2207 make the preconditioned conjugate gradient method resilient against multiple node failures.
2208 In *Proceedings of the 48th International Conference on Parallel Processing, ICPP 2019, Kyoto,*
2209 *Japan, August 05-08, 2019*, pages 67:1–67:10. ACM, 2019. doi:10.1145/3337821.3337849.
- 2210 **198** D. Stott Parker, Brad Pierce, and Paul R. Eggert. Monte Carlo arithmetic: a framework for
2211 the statistical analysis of roundoff error. *IEEE Computation in Science and Engineering*, 2001.
- 2212 **199** Alfredo Parra Hinojosa, Christoph Kowitz, Mario Heene, Dirk Pflüger, and H-J Bungartz.
2213 Towards a fault-tolerant, scalable implementation of GENE. In *Recent Trends in Computational*
2214 *Engineering-CE2014*, pages 47–65. Springer, 2015.
- 2215 **200** Sri Raj Paul, Akihiro Hayashi, Nicole Slattengren, Hemanth Kolla, Matthew Whitlock,
2216 Seonmyeong Bak, Keita Teranishi, Jackson Mayo, and Vivek Sankar. Enabling resilience in
2217 asynchronous many-task programming models. In Ramin Yahyapour, editor, *Euro-Par 2019:*
2218 *Parallel Processing*, pages 346–360. Springer International Publishing, 2019.
- 2219 **201** Benjamin Peherstorfer, Dirk Pflüger, and Hans-Joachim Bungartz. Density estimation
2220 with adaptive sparse grids for large data sets. In *Proceedings of the 2014 SIAM inter-*
2221 *national conference on data mining*, pages 443–451. SIAM, 2014. URL: [https://epubs.siam.](https://epubs.siam.org/doi/abs/10.1137/1.9781611973440.51)
2222 [org/doi/abs/10.1137/1.9781611973440.51](https://epubs.siam.org/doi/pdf/10.1137/1.9781611973440.51), arXiv:[https://epubs.siam.org/doi/pdf/10.](https://epubs.siam.org/doi/pdf/10.1137/1.9781611973440.51)
2223 [10.1137/1.9781611973440.51](https://epubs.siam.org/doi/pdf/10.1137/1.9781611973440.51), doi:10.1137/1.9781611973440.51.
- 2224 **202** A. J. Pena, W. Bland, and P. Balaji. VOCL-FT: Introducing techniques for efficient soft error
2225 coprocessor recovery. In *Proceedings of the International Conference on High Performance*
2226 *Computing, Networking, Storage and Analysis (SC'15)*, pages 1–12, Nov 2015. doi:10.1145/
2227 2807591.2807640.

- 2228 **203** M.D. Piggott, P.E. Farrell, C.R. Wilson, G.J. Gorman, and C.C. Pain. Anisotropic mesh
2229 adaptivity for multi-scale ocean modelling. *Philosophical Transactions of the Royal Society A:*
2230 *Mathematical, Physical and Engineering Sciences*, 367(1907):4591–4611, 2009.
- 2231 **204** A. Poulos, D. Wallace, R. Robey, L. Monroe, V. Job, S. Blanchard, W. Jones, and N. De-
2232 Bardeleben. Improving application resilience by extending error correction with contextual
2233 information. In *2018 IEEE/ACM 8th Workshop on Fault Tolerance for HPC at eXtreme Scale*
2234 *(FTXS)*, pages 19–28, Nov 2018. doi:10.1109/FTXS.2018.00006.
- 2235 **205** Karl E. Prikopa and Wilfried N. Gansterer. Fault-tolerant least squares solvers for wireless
2236 sensor networks based on gossiping. *J. Parallel Distributed Comput.*, 136:52–62, 2020. doi:
2237 10.1016/j.jpdc.2019.09.006.
- 2238 **206** Alfio Quarteroni and Alberto Valli. *Numerical approximation of partial differential equations*,
2239 volume 23. Springer Science & Business Media, 2008.
- 2240 **207** Petar Radojkovic, Manolis Marazakis, Paul Carpenter, Reiley Jeyapaul, Dimitris Gizopoulos,
2241 Martin Schulz, Adria Armejach, Eduard A Ayguade, François Bodin, Ramon Canal, Franck
2242 Cappello, Fabien Chaix, Guillaume Colin De Verdiere, Said Derradji, Stefano Di Carlo,
2243 Christian Engelmann, Ignacio Laguna, Miquel Moreto, Onur Mutlu, Lazaros Papadopoulos,
2244 Olly Perks, Manolis Ploumidis, Bezhad Salami, Yanos Sazeides, Dimitrios Soudris, Yiannis
2245 Soudris, Per Stenstrom, Samuel Thibault, Will Toms, and Osman Unsal. Towards Resilient
2246 EU HPC Systems: A Blueprint. Research report, European HPC resilience initiative, April
2247 2020. URL: <https://hal.inria.fr/hal-02922257>.
- 2248 **208** M. Wasiur Rashid and Michael C. Huang. Supporting highly-decoupled thread-level redundancy
2249 for parallel programs. In *Conf. on High-Performance Computer Architecture (HPCA)*, pages
2250 393–404. IEEE, 2008.
- 2251 **209** G. A. Reis, J. Chang, N. Vachharajani, R. Rangan, and D. I. August. Swift: software imple-
2252 mented fault tolerance. In *International Symposium on Code Generation and Optimization*,
2253 pages 243–254, 2005.
- 2254 **210** F. Rizzi, K. Morris, K. Sargsyan, P. Mycek, C. Safta, O. Le Maître, O.M. Knio, and B.J. De-
2255 busschere. Exploring the interplay of resilience and energy consumption for a task-based partial
2256 differential equations preconditioner. *Parallel Computing*, 73:16 – 27, 2018. Parallel Program-
2257 ming for Resilience and Energy Efficiency. URL: [http://www.sciencedirect.com/science/](http://www.sciencedirect.com/science/article/pii/S0167819117300753)
2258 [article/pii/S0167819117300753](http://www.sciencedirect.com/science/article/pii/S0167819117300753), doi:<https://doi.org/10.1016/j.parco.2017.05.005>.
- 2259 **211** Francesco Rizzi, Karla Morris, Khachik Sargsyan, Paul Mycek, Cosmin Safta, Bert Debusschere,
2260 Olivier LeMaitre, and Omar Knio. ULFM-MPI implementation of a resilient task-based partial
2261 differential equations preconditioner. In *Proceedings of the ACM Workshop on Fault-Tolerance*
2262 *for HPC at Extreme Scale*, New York, NY, USA, 2016. Association for Computing Machinery.
2263 doi:10.1145/2909428.2909429.
- 2264 **212** Gabriel Rodríguez, María J. Martín, Patricia González, Juan Touriño, and Ramón Doallo.
2265 Cppc: A compiler-assisted tool for portable checkpointing of message-passing applications.
2266 *Concurr. Comput.: Pract. Exper.*, 22(6):749–766, April 2010.
- 2267 **213** T. Ropars, T. V. Martsinkevich, A. Guermouche, A. Schiper, and F. Cappello. SPBC:
2268 Leveraging the characteristics of MPI HPC applications for scalable checkpointing. In *SC '13:*
2269 *Proceedings of the International Conference on High Performance Computing, Networking,*
2270 *Storage and Analysis*, pages 1–12, Nov 2013. doi:10.1145/2503210.2503271.
- 2271 **214** Ulrich Rüde. Fully adaptive multigrid methods. *SIAM Journal on Numerical Analysis*,
2272 30(1):230–248, 1993.
- 2273 **215** Ulrich Rüde. *Mathematical and computational techniques for multilevel adaptive methods*.
2274 SIAM, 1993.
- 2275 **216** Ulrich Rüde. Error estimators based on stable splittings. In David E Keyes and Jinchao Xu,
2276 editors, *Domain Decomposition Methods in Scientific and Engineering Computing: Proceedings*
2277 *of the Seventh International Conference on Domain Decomposition, October 27-30, 1993, the*
2278 *Pennsylvania State University*, volume 180 of *Contemporary Mathematics*, pages 111–118.
2279 American Mathematical Soc., 1994.

- 2280 217 Adrian Sandu. A class of multirate infinitesimal gark methods. *SIAM Journal on Numerical*
2281 *Analysis*, 57(5):2300–2327, 2019.
- 2282 218 Piyush Sao and Richard Vuduc. Self-stabilizing iterative solvers. *Proceedings of the*
2283 *Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems - ScalA '13*,
2284 pages 1–8, 2013. URL: <http://dl.acm.org/citation.cfm?doid=2530268.2530272>, doi:
2285 10.1145/2530268.2530272.
- 2286 219 Kento Sato, Naoya Maruyama, Kathryn Mohror, Adam Moody, Todd Gamblin, Bronis R
2287 de Supinski, and Satoshi Matsuoka. Design and modeling of a non-blocking checkpointing
2288 system. In *SC'12: Proceedings of the International Conference on High Performance Computing,*
2289 *Networking, Storage and Analysis*, pages 1–10. IEEE, 2012.
- 2290 220 V. Savcenco, W. Hundsdorfer, and J.G. Verwer. A multirate time stepping strategy for stiff
2291 ordinary differential equations. *BIT*, 47:137–155, 2007.
- 2292 221 Florian Schornbaum and Ulrich Rude. Massively parallel algorithms for the lattice Boltzmann
2293 method on nonuniform grids. *SIAM Journal on Scientific Computing*, 38(2):C96–C126, 2016.
- 2294 222 Florian Schornbaum and Ulrich Rude. Extreme-scale block-structured adaptive mesh refine-
2295 ment. *SIAM Journal on Scientific Computing*, 40(3):C358–C387, 2018.
- 2296 223 M. Schölzel. Reduced triple modular redundancy for built-in self-repair in vliw-processors. In
2297 *Signal Processing Algorithms, Architectures, Arrangements, and Applications SPA 2007*, pages
2298 21–26, Sep. 2007. doi:10.1109/SPA.2007.5903294.
- 2299 224 Bruno Seny, Jonathan Lambrechts, Thomas Toulorge, Vincent Legat, and Jean-François
2300 Remacle. An efficient parallel implementation of explicit multirate Runge–Kutta schemes for
2301 discontinuous Galerkin computations. *Journal of Computational Physics*, 256:135–160, 2014.
- 2302 225 F. Shahzad, J. Thies, M. Kreutzer, T. Zeiser, G. Hager, and G. Wellein. CRAFT: A library for
2303 easier application-level checkpoint/restart and automatic fault tolerance. *IEEE Transactions*
2304 *on Parallel and Distributed Systems*, 30(3):501–514, March 2019. doi:10.1109/TPDS.2018.
2305 2866794.
- 2306 226 Faisal Shahzad, Jonas Thies, Moritz Kreutzer, Thomas Zeiser, Georg Hager, and Gerhard
2307 Wellein. CRAFT: A library for easier application-level checkpoint/restart and automatic
2308 fault tolerance. *IEEE Transactions on Parallel and Distributed Systems*, 30(3):501–514, 2018.
2309 doi:10.1109/TPDS.2018.2866794.
- 2310 227 Manu Shantharam, Sowmyalatha Srinivasmurthy, and Padma Raghavan. Characterizing
2311 the impact of soft errors on iterative methods in scientific computing. *Proceedings of the*
2312 *international conference on Supercomputing - ICS '11*, page 152, 2011. URL: [http://portal.](http://portal.acm.org/citation.cfm?doid=1995896.1995922)
2313 [acm.org/citation.cfm?doid=1995896.1995922](http://portal.acm.org/citation.cfm?doid=1995896.1995922), doi:10.1145/1995896.1995922.
- 2314 228 Richard V. Southwell. *Relaxation methods in engineering science, a treatise on approximate*
2315 *computation*. Oxford, Oxford Univ. Pr., 1946. 1.ed., reprint.
- 2316 229 Robert Speck and Daniel Ruprecht. Toward fault-tolerant parallel-in-time integration with
2317 PFASST. *Parallel Computing*, 62:20–37, 2017.
- 2318 230 P. Spitéri. Parallel asynchronous algorithms for solving boundary value problems. In M. Cosnard
2319 et al., editor, *Parallel Algorithms*, pages 73–84. North-Holland, 1986.
- 2320 231 P. Spiteri. Parallel asynchronous algorithms: A survey. *Advances in Engineering Software*,
2321 149:1 – 34, 2020.
- 2322 232 Linda Stals. *Parallel multigrid on unstructured grids using adaptive finite element methods*.
2323 PhD thesis, Department Of Mathematics, The Australian National University, Australia, 1995.
- 2324 233 Linda Stals. Algorithm-based fault recovery of adaptively refined parallel multigrid grids. *The*
2325 *International Journal of High Performance Computing Applications*, 33(1):189–211, 2019.
- 2326 234 J. Stearley, K. Ferreira, D. Robinson, J. Laros, K. Pedretti, D. Arnold, P. Bridges, and
2327 R. Riesen. Does partial replication pay off? In *IEEE/IFIP International Conference on*
2328 *Dependable Systems and Networks Workshops (DSN 2012)*, pages 1–6, June 2012. doi:
2329 10.1109/DSNW.2012.6264669.
- 2330 235 Hana Straková, Gerhard Niederbrucker, and Wilfried N. Gansterer. Fault tolerance properties
2331 of gossip-based distributed orthogonal iteration methods. In Vassil N. Alexandrov, Michael Lees,

- 2332 Valeria V. Krzhizhanovskaya, Jack J. Dongarra, and Peter M. A. Sloot, editors, *Proceedings*
2333 *of the International Conference on Computational Science, ICCS 2013, Barcelona, Spain,*
2334 *5-7 June, 2013*, volume 18 of *Procedia Computer Science*, pages 189–198. Elsevier, 2013.
2335 doi:10.1016/j.procs.2013.05.182.
- 2336 **236** O. Subasi, J. Arias, O. Unsal, J. Labarta, and A. Cristal. Nan checkpoints: A task-based
2337 asynchronous dataflow framework for efficient and scalable checkpoint/restart. In *2015 23rd*
2338 *Euromicro International Conference on Parallel, Distributed, and Network-Based Processing*,
2339 pages 99–102, 2015.
- 2340 **237** O. Subasi, T. Martsinkevich, F. Zyulkyarov, O. Unsal, J. Labarta, and F. Cappello. Uni-
2341 fied fault-tolerance framework for hybrid task-parallel message-passing applications. *The*
2342 *International Journal of High Performance Computing Applications*, 32(5):641–657, 2018.
- 2343 **238** O. Subasi, G. Yalcin, F. Zyulkyarov, O. Unsal, and J. Labarta. A runtime heuristic to selectively
2344 replicate tasks for application-specific reliability targets. In *2016 IEEE International Conference*
2345 *on Cluster Computing (CLUSTER)*, pages 498–505, 2016.
- 2346 **239** O. Subasi, G. Yalcin, F. Zyulkyarov, O. Unsal, and J. Labarta. Designing and modelling
2347 selective replication for fault-tolerant HPC applications. In *17th IEEE/ACM International*
2348 *Symposium on Cluster, Cloud and Grid Computing (CCGRID)*, pages 452–457, May 2017.
2349 doi:10.1109/CCGRID.2017.40.
- 2350 **240** Omer Subasi, Javier Arias, Osman Unsal, Jesus Labarta, and Adrian Cristal. Programmer-
2351 directed partial redundancy for resilient HPC. In *Proceedings of the 12th ACM International*
2352 *Conference on Computing Frontiers (CF)*, pages 47:1–47:2, New York, NY, USA, 2015. ACM.
2353 URL: <http://doi.acm.org/10.1145/2742854.2742903>, doi:10.1145/2742854.2742903.
- 2354 **241** N. Sukhija, I. Banicescu, and F. M. Ciorba. Investigating the resilience of dynamic loop
2355 scheduling in heterogeneous computing systems. In *Proceedings of the 14th International*
2356 *Symposium on Parallel and Distributed Computing (ISPDC 2015)*, pages 194–203, June 2015.
2357 doi:10.1109/ISPDC.2015.29.
- 2358 **242** Daniel B. Szyld. The mystery of asynchronous iterations convergence when the spectral radius
2359 is one. Technical Report 98-102, Department of Mathematics, Temple University, Philadelphia,
2360 Pa., October 1998. Available at <http://www.math.temple.edu/szyld>.
- 2361 **243** Daniel B. Szyld. Perspectives on asynchronous computations for fluid flow problems. In K. J.
2362 Bathe, editor, *Computational Fluid and Solid Mechanics*, pages 377–380. Elsevier, 2001.
- 2363 **244** Daniel B. Szyld and Jian-Jun Xu. Convergence of some asynchronous nonlinear multisplitting
2364 methods. *Numerical Algorithms*, 25:347–361, 2000.
- 2365 **245** X. Tang, J. Zhai, B. Yu, W. Chen, W. Zheng, and K. Li. An efficient in-memory checkpoint
2366 method and its practice on fault-tolerant HPL. *IEEE Transactions on Parallel and Distributed*
2367 *Systems*, 29(4):758–771, April 2018. doi:10.1109/TPDS.2017.2781257.
- 2368 **246** D. Tao, S. Di, Z. Chen, and F. Cappello. Significantly improving lossy compression for scientific
2369 data sets based on multidimensional prediction and error-controlled quantization. In *2017*
2370 *IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 1129–1139.
2371 IEEE, 2017.
- 2372 **247** Dingwen Tao, Shuaiwen Leon Song, Sriram Krishnamoorthy, Panruo Wu, Xin Liang, Eddy Z.
2373 Zhang, Darren J. Kerbyson, and Zizhong Chen. New-sum: A novel online ABFT scheme for
2374 general iterative methods. In Hiroshi Nakashima, Kenjiro Taura, and Jack Lange, editors,
2375 *Proceedings of the 25th ACM International Symposium on High-Performance Parallel and*
2376 *Distributed Computing, HPDC 2016, Kyoto, Japan, May 31 - June 04, 2016*, pages 43–55. ACM,
2377 2016. iterative methods for linear systems: Krylov, stationary. doi:10.1145/2907294.2907306.
- 2378 **248** Containment Domain Team. Containment domains, 2014. URL: [https://lph.ece.utexas.](https://lph.ece.utexas.edu/public/CDs/ContainmentDomains)
2379 [edu/public/CDs/ContainmentDomains](https://lph.ece.utexas.edu/public/CDs/ContainmentDomains).
- 2380 **249** Nanos team. Nanos++, 2020. URL: [https://www.bsc.es/research-and-development/](https://www.bsc.es/research-and-development/software-and-apps/software-list/nanos-rtl)
2381 [software-and-apps/software-list/nanos-rtl](https://www.bsc.es/research-and-development/software-and-apps/software-list/nanos-rtl).
- 2382 **250** RAJA team. Raja performance portability layer. <https://github.com/LLNL/RAJA>, 2019.

- 2383 251 The Mercury Team. MERCURY programming language, 2020. URL: <https://www.mercurylang.org/>.
2384
- 2385 252 The Zoltan Team. Zoltan: Parallel Partitioning, Load Balancing and Data-Management
2386 Services, 2013. URL: <https://cs.sandia.gov/Zoltan/>.
- 2387 253 Keita Teranishi and Michael A. Heroux. Toward local failure local recovery resilience model
2388 using MPI-ULFM. In *Proceedings of the 21st European MPI Users' Group Meeting, EuroMPI/ASIA'14*,
2389 page 51–56, New York, NY, USA, 2014. Association for Computing Machinery. doi:10.1145/2642769.2642774.
2390
- 2391 254 Samuel Thibault. *On Runtime Systems for Task-based Programming on Heterogeneous Plat-*
2392 *forms*. Habilitation à diriger des recherches, Université de Bordeaux, December 2018. URL:
2393 <https://hal.inria.fr/tel-01959127>.
- 2394 255 G. Tumolo and L. Bonaventura. A semi-implicit, semi-Lagrangian discontinuous Galerkin
2395 framework for adaptive numerical weather prediction. *Quarterly Journal of the Royal Meteorological Society*,
2396 141(692):2582–2601, 2015.
- 2397 256 Henk A. van der Vorst and Qiang Ye. Residual replacement strategies for Krylov subspace
2398 iterative methods for the convergence of true residuals. *SIAM Journal on Scientific Computing*,
2399 22(3):835–852, 2000.
- 2400 257 Jyothish Varma, Chao Wang, Frank Mueller, Christian Engelmann, and Stephen L. Scott.
2401 Scalable, fault-tolerant membership for MPI tasks on HPC systems. In *Proceedings of the*
2402 *20th ACM International Conference on Supercomputing (ICS) 2006*, pages 219–228, Cairns,
2403 Australia, June 28-30, 2006. ACM Press, New York, NY, USA. doi:10.1145/1183401.
2404 1183433.
- 2405 258 J. Vignes. Discrete Stochastic Arithmetic for validating results of numerical software. *Numerical*
2406 *Algorithms*, 37(1–4):377–390, December 2004.
- 2407 259 T. N. Vijaykumar, Irith Pomeranz, and Karl Cheng. Transient-fault recovery using simultaneous
2408 multithreading. In *In proceedings of the 29th annual international symposium on computer*
2409 *architecture*, pages 87–98. IEEE Computer Society, 2002.
- 2410 260 John von Neumann. First Draft of a Report on the EDVAC, 1945. URL:
2411 [https://nsu.ru/xmlui/bitstream/handle/nsu/9018/2003-08-TheFirstDraft.pdf?](https://nsu.ru/xmlui/bitstream/handle/nsu/9018/2003-08-TheFirstDraft.pdf?sequence=1&isAllowed=y)
2412 [sequence=1&isAllowed=y](https://nsu.ru/xmlui/bitstream/handle/nsu/9018/2003-08-TheFirstDraft.pdf?sequence=1&isAllowed=y).
- 2413 261 John von Neumann. Probabilistic logics and the synthesis of reliable organisms from unreliable
2414 components. *Automata Studies*, pages 43–98, 1956.
- 2415 262 Chao Wang, Frank Mueller, Christian Engelmann, and Stephen L. Scott. A job pause
2416 service under LAM/MPI+BLCR for transparent fault tolerance. In *Proceedings of the 21st*
2417 *IEEE International Parallel and Distributed Processing Symposium (IPDPS) 2007*, pages
2418 1–10, Long Beach, CA, USA, March 26-30, 2007. ACM Press, New York, NY, USA. doi:
2419 10.1109/IPDPS.2007.370307.
- 2420 263 Chao Wang, Frank Mueller, Christian Engelmann, and Stephen L. Scott. Hybrid check-
2421 pointing for MPI jobs in HPC environments. In *Proceedings of the 16th IEEE Interna-*
2422 *tional Conference on Parallel and Distributed Systems (ICPADS) 2010*, pages 524–533,
2423 Shanghai, China, December 8-10, 2010. IEEE Computer Society, Los Alamitos, CA, USA.
2424 doi:10.1109/ICPADS.2010.48.
- 2425 264 Chao Wang, Frank Mueller, Christian Engelmann, and Stephen L. Scott. Proactive process-level
2426 live migration and back migration in HPC environments. *Journal of Parallel and Distributed*
2427 *Computing (JPDC)*, 72(2):254–267, February 2012. doi:10.1016/j.jpdc.2011.10.009.
- 2428 265 Chris Weaver and Todd M. Austin. A fault tolerant approach to microprocessor design.
2429 In *Proceedings of the 2001 International Conference on Dependable Systems and Networks*
2430 *(Formerly: FTCS)*, USA, 2001. IEEE Computer Society.
- 2431 266 Jordi Wolfson-Pou and Edmond Chow. Asynchronous multigrid methods. *2019 IEEE*
2432 *International Parallel and Distributed Processing Symposium (IPDPS)*, pages 101–110, 2019.
- 2433 267 Panruo Wu and Zizhong Chen. FT-ScaLAPACK: Correcting soft errors on-line for ScaLAPACK
2434 Cholesky, QR, and LU factorization routines. In *Proceedings of the 23rd International*

- 2435 *Symposium on High-performance Parallel and Distributed Computing*, HPDC '14, pages 49–60,
2436 New York, NY, USA, 2014. ACM. URL: <http://doi.acm.org/10.1145/2600212.2600232>,
2437 doi:10.1145/2600212.2600232.
- 2438 **268** Panruo Wu, Qiang Guan, Nathan DeBardeleben, Sean Blanchard, Dingwen Tao, Xin Liang,
2439 Jieyang Chen, and Zizhong Chen. Towards practical algorithm based fault tolerance in dense
2440 linear algebra. In *Proceedings of the 25th ACM International Symposium on High-Performance*
2441 *Parallel and Distributed Computing*, New York, NY, USA, 2016. Association for Computing
2442 Machinery. dense linear algebra. doi:10.1145/2907294.2907315.
- 2443 **269** Y. Kim, J. S. Plank, and J. J. Dongarra. Fault tolerant matrix operations using checksum and
2444 reverse computation. In *Proceedings of 6th Symposium on the Frontiers of Massively Parallel*
2445 *Computation (Frontiers '96)*, pages 70–77, Oct 1996. doi:10.1109/FMPC.1996.558063.
- 2446 **270** Ichitaro Yamazaki, Edmond Chow, Aurelien Bouteiller, and Jack J. Dongarra. Performance of
2447 asynchronous optimized Schwarz with one-sided communication. *Parallel Comput.*, 86:66–81,
2448 2019. doi:10.1016/j.parco.2019.05.004.
- 2449 **271** Yonghong Yan, Jisheng Zhao, Yi Guo, and Vivek Sarkar. Hierarchical place trees: A portable
2450 abstraction for task parallelism and data movement. In *International Workshop on Languages*
2451 *and Compilers for Parallel Computing*, pages 172–187. Springer, 2009.
- 2452 **272** Ulrike Yang, Piotr Luszczek, Satish Baley, and Keita Teranishi. An introduction to the xsdk a
2453 community of diverse numerical hpc software packages. Technical report, Sandia National
2454 Lab.(SNL-CA), Livermore, CA (United States), 2019. doi:10.6084/m9.figshare.7779452.
2455 v1.
- 2456 **273** J. Yu, D. Jian, Z. Wu, and H. Liu. Thread-level redundancy fault tolerant CMP based on
2457 relaxed input replication. In *2011 6th International Conference on Computer Sciences and*
2458 *Convergence Information Technology (ICCIT)*, pages 544–549, Nov 2011.
- 2459 **274** Gengbin Zheng, Xiang Ni, and Laxmikant V Kalé. A scalable double in-memory checkpoint
2460 and restart scheme towards exascale. In *IEEE/IFIP International Conference on Dependable*
2461 *Systems and Networks Workshops (DSN 2012)*, pages 1–6. IEEE, 2012.
- 2462 **275** Ziming Zheng, Andrew A. Chien, and Keita Teranishi. Fault tolerance in an inner-outer
2463 solver: A GVR-enabled case study. In Michel J. Daydé, Osni Marques, and Kengo Nakajima,
2464 editors, *High Performance Computing for Computational Science - VECPAR 2014 - 11th*
2465 *International Conference, Eugene, OR, USA, June 30 - July 3, 2014, Revised Selected Papers*,
2466 volume 8969 of *Lecture Notes in Computer Science*, pages 124–132. Springer, 2014. doi:
2467 10.1007/978-3-319-17353-5_11.
- 2468 **276** G Zoutendijk. Nonlinear programming, computational methods. *Integer and nonlinear*
2469 *programming*, pages 37–86, 1970.