

# Toward Resilient Algorithms and Applications

Michael A. Heroux, Sandia National Laboratories

## I. INTRODUCTION

Since the early days of supercomputing<sup>1</sup>, large-scale computing platforms have been engineered to handle unreliability. In contrast, algorithms and applications for large-scale systems have generally assumed a fairly simplistic failure model: The computer is a reliable digital machine, with consistent execution times and infrequent failures. If failure occurs, recovery can be handled by checkpoint-restart (CPR): occasionally storing a snapshot of application state and restarting from that saved state.

Over the past decade, the high performance computing community has become increasingly concerned that preserving the reliable, digital machine model will become too costly or infeasible [1], [2]. With the push toward exascale computing, this concern has become even greater [3], and we must explore other models and improve algorithms.

In this paper we discuss possibilities for developing new algorithms that are resilient to hard and soft failures. However, in order to reason about such algorithms, we first need programming models that enable more sophisticated recovery strategies than CPR.

## II. FOUR RESILIENCE-ENABLING PROGRAMMING MODELS

Algorithm-based fault tolerance has certainly been studied, going back many decades [4], and many algorithms have been developed [5], [6], [7], [8], but none of these algorithms have made it into broad practical use because we have no standard programming model support. In order to develop effective resilient algorithms and applications, we first need programming models that permit us to reason about failure and implement recovery. Here we present four programming models that we think have strong promise of being useful, ordering them from easiest to hardest to implement in a production system. Even though these models are not widely deployed today, using them as abstractions for developing new algorithms will provide motivation and guidance for development of both new algorithms and the underlying system software and hardware.

Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration (NNSA) under contract DE-AC04-94AL85000.

<sup>1</sup>The first Cray-1 (SN1) was delivered to Los Alamos National Laboratory without SECCED memory. Errors were so frequent that SN2 was scrapped and SN3 was delivered to NCAR with SECCED.

### A. *Skeptical Programming (SkP)*

Almost all algorithm developers assume that their software will execute reliably or fail obviously by halting. However, if they understand that silent data corruption is possible, they can develop very simple and inexpensive validation tests based on their understanding of the mathematical properties of their algorithms. For example, many algorithms have global properties of orthogonality, or conservation principles, that are implicitly assumed to be true during the execution of an algorithm. If these properties were checked occasionally during execution, the cost can be very low and many silent data corruption events could be detected. Recovery may be as simple as aborting, or may involve rolling back to a previous valid state, or even continuing execution if the error will be damped by subsequent computations.

### B. *Relaxed Bulk-synchronous Programming (RBSP)*

One of the first impacts of reduced reliability is performance variability. As low-level system failure rates increase, error detection and correction happen more frequently in the hardware and system software layers. These events preserve the reliable digital machine model, but introduce variability in execution time. Many scalable applications are designed under the implicit assumption that equal work implies equal execution time, so that if we balance the work of a parallel application, we should scale well on a parallel computer, even though we synchronize across processors during execution. Performance variability, when coupled with frequent collective operations, leads to severe limitations in scalability, especially as we go to a million or more processes.

With the introduction of MPI-3 [9], we now have asynchronous neighborhood and global collectives, enabling a “relaxed” bulk-synchronous programming model (RBSP). Given RBSP capabilities, we are now able to develop new algorithms that can potentially hide latency.

### C. *Local Failure, Local Recovery (LFLR)*

For parallel applications based on MPI, the current approach to dealing with the loss of a single process is to kill all remaining processes and restart the application. As we now regularly run on hundreds of thousands to more than a million processors, this approach is not feasible. Instead a local failure should permit a local recovery<sup>2</sup>.

<sup>2</sup>By local recovery we do not mean that no communication is done in the detection or recovery phases. Instead we mean that recovery is logically local from an application developer's perspective. All processes that have valid state are not involved in recovery, except to the extent that they assist in restoring state on the failed process, or take over its workload.

One type of local-failure-local-recovery (LFLR) model permits the user to store specific data *persistently* for each MPI process and allows a recovery function to be registered, such that, if a process fails, a new process is started and assigned to the rank of the failed process, and the user's recovery function is called, giving access to the persistent data of the old process, as well as the neighbors' persistent data. Using LFLR, we can develop new algorithms for many types of problems.

#### D. Selective Reliability Programming (SRP)

The fourth programming model we discuss is Selective Reliability Programming (SRP), which gives the programmer the ability to declare specific data and compute regions to be more reliable than the "bulk" reliability of the underlying system (or we can switch the default to be reliable and then selectively be less reliable). By distinguishing between what needs to be highly reliable or not, we can develop new algorithms that store most data and do most computations with low reliability while retaining the robustness of a fully reliable approach.

Although the costs of high reliability will impact the practicality of some approaches, the details of how reliability is implemented is not fundamentally important to reasoning about new algorithms. In some cases, even very expensive approaches such as triple modular redundancy (TMR) can still be much faster than a fully unreliable approach.

### III. TOWARD RESILIENT ALGORITHMS

Resilient algorithms have long been a subject of research. The above four programming models enable further research and drive co-development of the algorithms and computing system features that are required to realize resilient applications. In this section we discuss some of the many possible algorithms that can be developed under the above programming models in order to provide resilience on future systems.

#### A. Detecting and Responding to Silent Data Corruption

Skeptical programming can be used to detect silent data corruption such as bit flips, and then determine if the resulting error is "harmless" or not. One example of this kind of algorithm is an implementation of GMRES [10] that detects and, optionally, corrects single bit flips very inexpensively as part of the Arnoldi process. Many existing ABFT algorithms can be implemented in using a skeptical algorithm programming approach, since the meta data used to recover state can also be used to detect anomalous behavior.

#### B. Latency-tolerant Algorithms

One of the most important and effective algorithm research and development strategies we can explore now is latency tolerance. Many of our scalable algorithms and applications depend on collective operations that, when implemented in a straightforward manner, lead to synchronous global

collectives. On emerging high end platforms, these collectives have become severe performance limiters due to poor scaling of collectives. The advent of asynchronous collectives gives us new opportunities. The basic challenge we face as algorithm developers in this situation is finding useful work to do while a collective is completing. Recent work in pipelining algorithms, for example the  $p(l)$ -GMRES algorithm [11], shows that latency hiding by unrolling iterations in a Krylov solver can help restore scalability. Similar approaches for many algorithms can lead to relatively minor design changes that result in better tolerance of latency and performance variability.

The impact of successfully redesigning algorithms to be latency tolerant is that performance variability on existing systems can be hidden. But even more importantly, if we can tolerate performance variability due to error detection and correction at the system software and hardware levels, system designers can detect and correct more errors without impacting application scalability, permitting us to extend the viability of the reliable digital machine model.

#### C. Locally Restarted PDE Computations

Given the programming features described in Section II-C, we have the potential to develop a broad collection of algorithm with local recovery properties. Examples for differential equations include:

- Explicit methods: As shown in [12], an explicit time-stepping algorithm can be easily implemented to recover locally, given the LFLR features.
- Implicit methods: This case is more interesting. The challenge is to restore a local state that is equivalent up to the truncation error of the PDE. Several interesting approaches seem promising.
- Redundant storage of coarse model: In order to recover state from a lost process, we could explore storing a coarse model representation on neighboring processes that could be used to boot-strap state recovery upon failure.

#### D. Reliable Outer Iterations

Many algorithms can be cast in an outer-inner formulation. For example, a fault-tolerant GMRES variant, as described in [13], uses reliable computation and storage in the outer iteration and an "unreliable" GMRES in the inner iteration. The result is that most computation and data are in low-reliability mode, leading to presumably cheaper computations. Because the outer iterations are reliable, the solution returned by the inner solve (if it comes back at all) can be analyzed and used or discarded. Even if the inner solve answer is not correct, it can still be used with some effect.

### IV. CONCLUSIONS

Resilience is a critical requirement for future high-end computing. In order to effectively develop resilient algorithms and applications, we need robust and usable

resilient computing models. In this paper we have identified four specific models that allow us to reason about and develop new algorithms. SkP requires nothing more than a change in attitude on the part of the programmer, from trusting that a machine is reliable digital device to being aware that an incorrect computation may occur. RBSP is already possible with the introduction of MPI 3.0. LFLR requires more support from the underlying system layers, and requires some kind of support from programming languages and libraries. The ULFM library [14], [15] already provides one approach to supporting LFLR. SRP is the most challenging model, but also firmly addresses one of the biggest challenges we face: silent errors.

Much of the focus of extreme-scale computing is on massive concurrency, which is appropriate. However, without resilient computing models we face a very real risk of application failure rates that are too high to realize the benefits of future systems. In addition, any progress we make in resilient algorithms and applications permits us to utilize lower cost systems in general, systems with lower quality interconnect networks, higher bit failure rates and increased node loss. Resilient algorithms and applications will enable effective extreme scale computing and reduce system costs at other levels.

## REFERENCES

- [1] N. Miskov-Zivanov and D. Marculescu, "Soft error rate analysis for sequential circuits," in *Proceedings of the Conference on Design, Automation and Test in Europe*, ser. DATE '07. San Jose, CA, USA: EDA Consortium, 2007, pp. 1436–1441. [Online]. Available: <http://portal.acm.org/citation.cfm?id=1266366.1266680>
- [2] T. Karnik, P. Hazucha, and J. Patel, "Characterization of soft errors caused by single event upsets in CMOS processes," *IEEE Trans. Dependable Secur. Comput.*, vol. 1, pp. 128–143, April 2004. [Online]. Available: <http://dx.doi.org/10.1109/TDSC.2004.14>
- [3] e. a. John Daly, "Inter-agency workshop on hpc resilience at extreme scale." [Online]. Available: <http://institute.lanl.gov/resilience/docs/Inter-AgencyResilienceReport.pdf>
- [4] K.-H. Huang and J. A. Abraham, "Algorithm-based fault tolerance for matrix operations," *IEEE Transactions on Computers*, vol. C-33, no. 6, June 1984.
- [5] G. A. Reis, J. Chang, N. Vachharajani, R. Rangan, and D. I. August, "Swift: Software implemented fault tolerance," in *Proceedings of the international symposium on Code generation and optimization*, ser. CGO '05. Washington, DC, USA: IEEE Computer Society, 2005, pp. 243–254. [Online]. Available: <http://dx.doi.org/10.1109/CGO.2005.34>
- [6] J. Langou, Z. Chen, G. Bosilca, and J. Dongarra, "Recovery patterns for iterative methods in a parallel unstable environment," *SIAM J. Sci. Comput.*, vol. 30, pp. 102–116, November 2007. [Online]. Available: <http://portal.acm.org/citation.cfm?id=1350656.1350657>
- [7] T. Wilfredo, "Software fault tolerance: A tutorial," Tech. Rep., 2000.
- [8] X. Yang, Y. Du, P. Wang, H. Fu, J. Jia, Z. Wang, and G. Suo, "The fault tolerant parallel algorithm: the parallel recomputing based failure recovery," in *Proceedings of the 16th International Conference on Parallel Architecture and Compilation Techniques*, ser. PACT '07. Washington, DC, USA: IEEE Computer Society, 2007, pp. 199–212. [Online]. Available: <http://dx.doi.org/10.1109/PACT.2007.73>
- [9] MPI Forum, "MPI: A Message-Passing Interface Standard. Version 3.0," 2012, available at: <http://www.mpi-forum.org> (Sep. 2012).
- [10] J. N. C. S. U. Elliott, M. F. S. N. L. Hoemmen, and Muel, "Making GMRES Resilient to Single Bit Flips," 2013. [Online]. Available: [http://sc13.supercomputing.org/sites/default/files/PostersArchive/tech\\_posters/post219s2-file2.pdf](http://sc13.supercomputing.org/sites/default/files/PostersArchive/tech_posters/post219s2-file2.pdf)
- [11] P. Ghysels, T. J. Ashby, K. Meerbergen, and W. Vanroose, "Hiding global communication latency in the GMRES algorithm on massively parallel machines," Submitted to *SIAM J. Sci. Comp.* Available online at <http://www.ua.ac.be/download.aspx?c=pieter.ghysels&n=92953&ct=84288&e=292475> [last accessed 29 April 2013].
- [12] D. Wong and M. Gokhale, "A memory-mapped approach to checkpointing."
- [13] P. G. Bridges, K. B. Ferreira, M. A. Heroux, and M. Hoemmen, "Fault-tolerant linear solvers via selective reliability," *ArXiv e-prints*, Jun. 2012.
- [14] W. Bland, A. Bouteiller, T. Herault, G. Bosilca, and J. Dongarra, "Post-failure recovery of MPI communication capability: Design and rationale," *International Journal of High Performance Computing Applications*, vol. 27, no. 3, pp. 244–254, Jun. 2013. [Online]. Available: <http://hpc.sagepub.com/content/27/3/244>
- [15] "User Level Failure Mitigation — ICL Fault Tolerance." [Online]. Available: <http://fault-tolerance.org/ulfm/>