

# An Evaluation of User-Level Failure Mitigation Support in MPI

Wesley Bland<sup>1</sup>, Aurelien Bouteiller<sup>1</sup>, Thomas Herault<sup>1</sup>, Joshua Hursey<sup>2</sup>,  
George Bosilca<sup>1</sup>, and Jack J. Dongarra<sup>1</sup>

<sup>1</sup> Innovative Computing Laboratory, University of Tennessee  
{bland, bouteill, herault, bosilca, dongarra}@eecs.utk.edu

<sup>2</sup> Oak Ridge National Laboratory  
hurseyjj@ornl.gov

**Abstract.** As the scale of computing platforms becomes increasingly extreme, the requirements for application fault tolerance are increasing as well. Techniques to address this problem by improving the resilience of algorithms have been developed, but they currently receive no support from the programming model, and without such support, they are bound to fail. This paper discusses the failure-free overhead and recovery impact aspects of the User-Level Failure Mitigation proposal presented in the MPI Forum. Experiments demonstrate that fault-aware MPI has little or no impact on performance for a range of applications, and produces satisfactory recovery times when there are failures.

## 1 Introduction

In a constant effort to deliver steady performance improvements, the size of High Performance Computing (HPC) systems, as observed by the Top 500 ranking<sup>3</sup>, has grown tremendously over the last decade. This trend is unlikely to stop, as outlined by the International Exascale Software Project (IESP) [9] projection of the Exaflop platform, a milestone that should be reached as soon as 2019. Based on the foreseeable limits of the infrastructure costs, an Exaflop capable machine is expected to be built from gigahertz processing cores, with thousands of cores per computing node, thus requiring millions of computing cores to reach the mark. Even under the most optimistic assumptions about the individual components' reliability, probabilistic amplification from using millions of nodes has a dramatic impact on the Mean Time Between Failure (MTBF) of the entire platform. The probability of a failure happening *during the next hour* on an Exascale platform is disturbingly close to 1; thereby many computing nodes will inevitably fail during the execution of an application [7]. It is even more alarming that most popular fault tolerant approaches see their efficiency plummet at Exascale [3, 4], calling for application centric failure mitigation strategies [15].

The prevalence of distributed memory machines promotes the use of the message passing model. An extensive and varied spectrum of domain science ap-

---

<sup>3</sup> <http://www.top500.org/>

plications depend on libraries compliant with the MPI standard<sup>4</sup>. Although unconventional programming paradigms are emerging [18, 20], most delegate their data movements to MPI and it is widely acknowledged that MPI is here to stay. However, MPI has to evolve to effectively support the demanding requirements imposed by novel architectures, programming approaches, and dynamic runtime systems. In particular, its support for fault tolerance has always been inadequate [13]. To address the growing interest in fault-aware MPI, a working group has been formed in the context of the MPI Forum. Their User-Level Failure Mitigation (ULFM) [1] proposal features the basic interface and new semantics to enable applications and libraries to repair the state of MPI and tolerate failures. The purpose of this paper is to evaluate the tradeoffs that are needed for the integration of this fault mitigation specification and its impact (or lack thereof) on MPI performance and scalability. The contributions of this work are to evaluate the difficulties faced by MPI implementors, and demonstrate the feasibility of a low-impact implementation on the failure-free performance as well as an estimate of the recovery time of the MPI state after a failure.

The remainder of this paper is organized as follows: the next section introduces a short history of fault tolerance in MPI; Section 3 presents the constructs introduced by the proposal; Section 4 discusses the challenges faced by MPI implementors; then the performance impact of the implementation in Open MPI is discussed in Section 5 before we conclude in Section 6.

## 2 Related Work

Efforts toward fault tolerance in MPI have previously been attempted. Automatic fault tolerance [5, 6] is a compelling approach for users, as failures are completely masked and handled internally by the MPI library, which requires no new interfaces to MPI or application code changes. Unfortunately, many recent studies point out that automatic approaches, either based on checkpoints or replication, will exhibit poor efficiency on Exaflop platforms [3, 4].

Application Based Fault Tolerance (ABFT) [8, 10, 15] is another approach that promises better scalability, at the cost of significant algorithm and application code changes. Despite some limited successes [2, 13], MPI interfaces need to be extended to effectively support ABFT. The most notable past effort is FT-MPI [11]. Several recovery modes were available to the user. In the *Blank* mode, failed processes were replaced by `MPI_PROC_NULL`; messages to and from them were silently ignored and collective algorithms had to be significantly modified. In the *Replace* mode, faulty processes were replaced with new processes. In all cases, only `MPI_COMM_WORLD` would be repaired and the application was in charge of rebuilding any other communicators, leading to difficult library composition. No standardization effort was pursued, and it was mostly used as a playground for understanding the fundamental concepts.

A more recent effort to introduce failure handling mechanisms was the Run-Through Stabilization proposal [16]. This proposal introduced many new con-

---

<sup>4</sup> <http://mpi-forum.org/docs/mpi-2.2/mpi22-report.pdf>

structs for MPI including the ability to “validate” communicators as a way of marking failure as recognized and allowing the application to continue using the communicator. It included other new ideas such as Failure Handlers for uniform failure notification. Because of the implementation complexity imposed by resuming operations on failed communicators, this proposal was eventually unsuccessful in its introduction to the MPI Standard.

### 3 New MPI Constructs

This section succinctly presents the prominent interfaces proposed to enable effective support of User-Level Failure Mitigation for MPI applications. The interested reader can refer to the technical document for a complete description of the interfaces [1] and to the amended standard draft<sup>5</sup>.

Designing the mechanism that users would use to manage failures was built around three concepts: 1) simplicity, the API should be easy to understand and use in most common scenarios; 2) flexibility, the API should allow varied fault tolerant models to be built as external libraries and; 3) absence of deadlock, no MPI call (point-to-point or collective) can block indefinitely after a failure, but must either succeed or raise an MPI error. Two major pitfalls must be avoided: jitter prone, permanent monitoring of the health of peers a process is not actively communicating with, and expensive consensus required for returning consistent errors at all ranks. The operative principle is then that errors (`MPI_ERR_PROC_FAILED`) are not indicative of the return status on remote processes, but are raised only at a particular rank, when a particular operation cannot complete because a participating peer has failed. The following functions provide the basic blocks for maintaining consistency and enabling recovery of the state of MPI.

`MPI_COMM_FAILURE_ACK` & `MPI_COMM_FAILURE_GET_ACKED`: These two calls allow the application to determine which processes within a communicator have failed. The acknowledgement function serves to mark a point in time which will be used as a reference. The function to get the acknowledged failures refers back to this reference point and returns the group of processes which were locally known to have failed. After acknowledging failures, the application can resume `MPI_ANY_SOURCE` point-to-point operations between non-failed processes, but operations involving failed processes (such as collective operations) will likely continue to raise errors.

`MPI_COMM_REVOKE`: Because failure detection is not global to the communicator, some processes may raise an error for an operation, while others do not. This inconsistency in error reporting may result in some processes continuing their normal, failure-free execution path, while others have diverged to the recovery execution path. As an example, if a process, unaware of the failure, posts a reception from another process that has switched to the recovery path, the matching send will never be posted. Yet no failed process participates in the operation

---

<sup>5</sup> <http://svn.mpi-forum.org/trac/mpi-forum-web/ticket/323>

and it should not raise an error. The receive operation is effectively deadlocked. The revoke operation provides a mechanism for the application to resolve such situations before entering the recovery path. A revoked communicator becomes improper for further communication, and all future or pending communications on this communicator will be interrupted and completed with the new error code `MPI_ERR_REVOKED`. It is notable that although this operation is not collective (a process can enter it alone), it affects remote ranks without a matching call.

`MPI_COMM_SHRINK`: The shrink operation allows the application to create a new communicator by eliminating all failed processes from a revoked communicator. The operation is collective and performs a consensus algorithm to ensure that all participating processes complete the operation with equivalent groups in the new communicator. This function cannot return an error due to process failure. Instead, such errors are absorbed as part of the consensus algorithms and will be excluded from the resulting communicator.

`MPI_COMM_AGREE`: This operation provides an agreement algorithm which can be used to determine a consistent state between processes when such strong consistency is necessary. The function is collective and forms an agreement over a boolean value, even when failures have happened or the communicator has been revoked. The agreement can be used to resolve a number of consistency issues after a failure, such as uniform completion of an algorithmic phase or collective operation, or as a key building block for strongly consistent failure handling approaches (such as transactions).

## 4 Implementation Issues

In this section, we detail the challenges and advantages of the aforementioned MPI constructs. They unfold along three main axes, the amount of supplementary state and memory to be kept within the MPI library, the additional operations to be executed on the critical path of communication routines, and the algorithmic cost of failure recovery routines. We discuss, in general, options available to implementors, and highlight issues with insight from a prototype implementation in Open MPI [12].

### 4.1 Impact on communication routines

*Memory*: Because a communicator cannot be repaired, tracking the state of failed processes imposes a minimal memory overhead. From a practical perspective each node needs a global list of detected failures, shared by all communicators; its size grows linearly with the number of failures, and it is empty as long as no failures occur. Within each communicator, the supplementary state is limited to two values: whether the communicator is revoked or not, and an index in the global list of failures denoting the last acknowledged failure (with `MPI_COMM_FAILURE_ACK`). For efficiency reasons, an implementation may decide

to cache the fact that some failures have happened in the communicator so that collective operations and `MPI_ANY_SOURCE` receptions can bail out quickly. Overall, the supplementary memory consumption from fault tolerant constructs is small, independent of the total number of nodes, and unlikely to affect the cache and TLB hit rates.

*Conditionals:* Another concern is the number of supplementary conditions on the latency critical path. Indeed, most completion operations require a supplementary conditional statement to handle the case where the underlying communication context has been revoked. However, the prediction branching logic of the processor can be hinted to favor the failure free outcome, resulting in a single load of a cached value and a single, mostly well-predicted, branching instruction, unlikely to affect the instruction pipeline. It is notable that non-blocking operations raise errors related to process failure only during the completion step, and thus do not need to check for revocation before the latency critical section.

*Matching logic:* `MPI_COMM_REVOKE` does not have a matching call on other processes on which it has an effect. As such, it might add detrimental complexity to the matching logic. However, any MPI implementation needs to handle unexpected messages. The order of revocation message delivery is loose enough that the handling of revocation notices can be integrated within the existing unexpected message matching logic. In our implementation in Open MPI, we leverage the active message low level transport layer to introduce revocation as a new active message tag, without a single change to the matching logic.

*Collective operations:* A typical MPI implementation supports a large number of collective algorithms, which are dynamically selected depending on criteria such as communicator or message size and hardware topology. The loose requirements of the proposal concerning error reporting of process failures in collective operations limits the impact it has on collective operations. Typically, the collective communication algorithms and selection logic are left unchanged. The only new requirement is that failures happening at any rank of the communicator cause all processes to exit the collective (successfully for some, with an error for others). Due to the underlying loosely-connected topologies used by some algorithms, a point-to-point based implementation of a collective communication is unlikely to detect all process failures. Fortunately, a practical implementation exists that does not require modifying any of the collective operations: when a rank raises an error because of a process failure, it can revoke an internal, temporary communication context associated with the collective operation. As the revocation notice propagates on the internal communicator, it interrupts the point-to-point operations of the collective. An error code is returned to the high level MPI wrapper, which in turn raises the appropriate error on the user's communicator.

## 4.2 Recovery routines

Some of the recovery routines described in Section 3 are unique in their ability to deliver a valid result despite the occurrence of failures. This specification of

correct behavior across failures calls for resilient, more complex algorithms. In most cases, these functions are intended to be called sparingly by users, only after actual failures have happened, as a means of recovering a consistent state across all processes. The remainder of this section describes the algorithms that can be used to deliver this specification and their cost.

*Agreement:* The agreement can be conceptualized as a failure-resilient reduction on a boolean value. Many agreement algorithms have been proposed in the literature; the log-scaling two-phase consensus algorithm used by the ULFM prototype is one of many possible implementations of `MPI_COMM_AGREE` operation based upon prior work in the field. Specifically, this algorithm is a variation of the multi-level two-phase commit algorithms [19]. The algorithm first performs a reduction of the input values to an elected coordinator in the communicator. The coordinator then makes a decision on the output value and broadcasts that value back to all of the alive processes in the communicator. The complexity of the agreement algorithm appears when adapting to an emerging process failure of the coordinator and/or participants. A more extensive discussion of the algorithmic complexity has been published by Hursey, et.al. [17]. The algorithmic complexity of this implementation is  $O(\log(n))$  for the failure free case, matching that of an `MPI_ALLREDUCE` operation over the alive processes in the communicator.

*Revoke:* Although the revoke operation is not collective, the revocation notification needs to be propagated to all alive processes in the specified communicator, even when new failures happen during the revoke propagation. These requirements are not without recalling those from the *reliable broadcast* [14]. Among the four defining qualities of a reliable broadcast (*Termination, Validity, Integrity, Agreement*), the termination and integrity criteria can be relaxed in the context of the revoke algorithm. If a failure during the Revoke algorithm kills the initiator as well as all the already notified processes, the Revoke notification is indeed lost, but the observed behavior, from the view of the application, is indiscernible from a failure at the initiator before the propagation started. As the algorithm still ensures agreement, there are no opportunities for inconsistent views.

In the ULFM implementation, we used a naive flooding algorithm for simplicity. The initiator marks the communicator as revoked and sends a Revoke message to every processes in the groups (local and remote) of the communicator. Upon reception of a revoke message, if the communicator is not already revoked, it is revoked and the process acts as a new initiator. Better algorithms exist, but even this naive approach provides reasonable performance (see Section 5) considering it is called only in response to an actual failure.

*Shrink:* The Shrink operation is, algorithmically, an agreement on which the consensus is done on the group of failed processes. Hence, the two operations have the same algorithmic complexity. Indeed, in the prototype implementation, `MPI_COMM_AGREE` and `MPI_COMM_SHRINK` share the same internal implementation of the agreement.

## 5 Performance Analysis

The following analysis used a prototype of the ULFM proposal based on the development trunk of Open MPI [12] (r26237). The test results presented were gathered from the Smoky system at Oak Ridge National Laboratory. Each node contains four quad-core 2.0 GHz AMD Opteron processors with 2 GB of memory per compute core. Compute nodes are connected with gigabit Ethernet and InfiniBand. Some shared-memory benchmarks were conducted on Romulus, a  $6 \times 8$ -core AMD Opteron 6180 SE with 256GB of memory (32GB per socket) at the University of Tennessee.

The NetPIPE benchmark (v3.7) was used to assess the 1-byte latency and bandwidth impact of the modifications necessary for the ULFM support in Open MPI. We compare the vanilla version of Open MPI (r26237) with the ULFM enabled version on Smoky. Table 1 highlights the fact that the differences in performance are well below the noise limit, and that the standard deviation is negligible proving the performance stability and lack of impact.

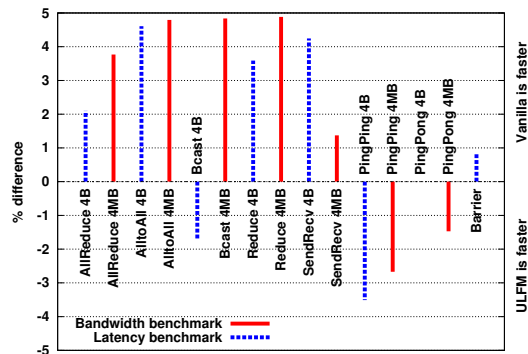
1-byte Latency (microseconds) (cache hot)					
Interconnect	Vanilla	Std. Dev.	Enabled	Std. Dev.	Difference
Shared Memory	0.8008	0.0093	0.8016	0.0161	0.0008
TCP	10.2564	0.0946	10.2776	0.1065	0.0212
OpenIB	4.9637	0.0018	4.9650	0.0022	0.0013

Bandwidth (Mbps) (cache hot)					
Interconnect	Vanilla	Std. Dev.	Enabled	Std. Dev.	Difference
Shared Memory	10,625.92	23.46	10,602.68	30.73	-23.24
TCP	6,311.38	14.42	6,302.75	10.72	-8.63
OpenIB	9,688.85	3.29	9,689.13	3.77	0.28

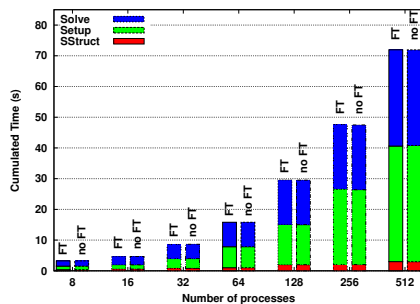
**Table 1.** NetPIPE results on Smoky.

The impact on shared memory systems, which are sensitive even to small modifications of the MPI library, has been further assessed on the Romulus machine – a large shared memory machine – using the IMB benchmark suite (v3.2.3). As shown in Figure 1, the duration difference of all the benchmarks (point-to-point and collective) remains below 5%, thus within the standard deviation of the implementation on that machine.

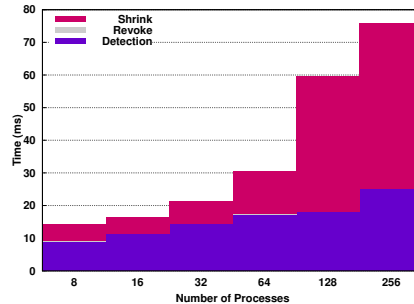


**Fig. 1.** The Intel MPI Benchmarks: relative difference between ULFM and the vanilla Open MPI on shared memory (Romulus). Standard deviation  $\approx 5\%$  on 1,000 runs.

To measure the impact of the prototype on a real application, we used the Sequoia AMG benchmark<sup>6</sup>. This MPI intensive benchmark is an Algebraic Multi-Grid (AMG) linear system solver for unstructured mesh physics. A weak scaling study was conducted up to 512 processes following the problem *Set 5*. In Figure 2, we compare the time slicing of three main phases (Solve, Setup, and SStruct) of the benchmark, with, side by side, the vanilla version of the Open MPI implementation, and the ULFM enabled one. The application itself is not fault tolerant and does not use the features proposed in ULFM. The goal of this benchmark is to demonstrate that a careful implementation of the proposed semantic does not impact the performance of the MPI implementation, and ultimately leaves the behavior and performance of legacy applications unchanged. The results show that the performance difference is negligible.



**Fig. 2.** Comparison of the vanilla and ULFM versions of Open MPI running Sequoia-AMG at different scales (Smoky).



**Fig. 3.** Evaluation of the Fault Injection Benchmark with full recovery at different scales (Smoky).

To assess the overheads of recovery constructs, we developed a synthetic benchmark that mimics the behavior of a typical fixed-size tightly-coupled fault-tolerant application. Unlike a normal application it performs an infinite loop, where each iteration contains a failure and the corresponding recovery procedure. Each iteration consists of 5 phases: in the first phase (*Detection*), all processes but a designated victim enter a Barrier on the intracommunicator. The victim dies, and the failure detection mechanism makes all surviving processes exit the Barrier, some with an error code. In Phase 2 (*Revoke*), the surviving processes that detected a process-failure related error during the previous phase invoke the new construct `MPI_COMM_REVOKE`. Then they proceed to Phase 3 (*Shrink*) where the intracommunicator is shrunk using `MPI_COMM_SHRINK`. The two other phases serve to repair a full-size intracommunicator using spawn and intercommunicator merge operations to allow the benchmark to proceed to the next round.

In Figure 3, we present the timing of each phase, averaged upon 50 iterations of the benchmark loop, for a varying number of processes on the Smoky machine. We focus on the three points related to ULFM: failure detection, revoke and shrink. The failure detection is mildly impacted by the scale. In the proto-

<sup>6</sup> <https://asc.llnl.gov/sequoia/benchmarks/#amg>



type implementation, the detection happens at two levels, either in the runtime system or in the MPI library (when it occurs on an active link). Between the two detectors, all ranks get notified within 30ms of the failure (this compares to the 1s timeout at the link level). Although the revoke call will inject a linear number of messages (at each rank) in the network to implement the level of reliability required for this operation, the duration of this call itself is under  $50\mu s$  and is not visible in the figure. The network is disturbed for a longer period, due to the processing of the messages, but this disturbance will appear in the network only after a failure occurred. The last call shown in the figure is the shrink operation. Although its duration increases linearly with the number of processes (the figure has a logarithmic scale on the x-axis), this cost must only be paid after a failure, in order to continue using collective operations. In its current implementation, shrink requires an agreement, the allocation of a new communicator identifier, and the creation of the communicator (with `MPI_COMM_SPLIT`). Most of the time spent in the shrink operation is not in the agreement (which scales logarithmically), but in the underlying implementation of the communicator creation.

## 6 Conclusion

Many responsible voices agree that sharp increases in the volatility of future, extreme scale computing platforms are likely to imperil our ability to use them for advanced applications that deliver meaningful scientific results and maximize research productivity. Moreover, it is clear that any techniques developed to address this volatility must be supported in the programming and execution model. Since MPI is currently, and will likely continue to be – in the medium-term – both the de-facto programming model for distributed applications and the default execution model for large scale platforms running at the bleeding edge, MPI is the place in the software infrastructure where semantic and run-time support for application faults needs to be provided.

The ULFM proposal is a careful but important step forward toward accomplishing this goal. It not only delivers support for a number of new and innovative resilience techniques, it provides this support through a simple, straightforward and familiar API that requires minimal modifications of the underlying MPI implementation. Moreover, it is backward compatible with previous versions of the MPI standard, so that non fault-tolerant applications (legacy or otherwise) are supported without any changes to the code. Perhaps most significantly, applications can use ULFM-enabled MPI without experiencing any degradation in their performance, as we demonstrate in this paper.

Several applications, ranging from Master-Worker to tightly coupled, are currently being refactored to take advantage of the semantics in this proposal. Beyond applications, the expressivity of this proposal is investigated in the context of providing extended fault tolerance models as convenience, portable libraries.

## References

1. Bland, W., Bosilca, G., Bouteiller, A., Herault, T., Dongarra, J.: A proposal for User-Level Failure Mitigation in the MPI-3 standard. Tech. rep., Department of Electrical Engineering and Computer Science, University of Tennessee (2012)
2. Bland, W., Du, P., Bouteiller, A., Herault, T., Bosilca, G., Dongarra, J.J.: A Checkpoint-on-Failure protocol for algorithm-based recovery in standard MPI. In: 18th Euro-Par. p. to appear. LNCS, Springer (2012)
3. Bosilca, G., Bouteiller, A., Brunet, É., Cappello, F., Dongarra, J., Guermouche, A., Héroult, T., Robert, Y., Vivien, F., Zaidouni, D.: Unified Model for Assessing Checkpointing Protocols at Extreme-Scale. Tech. report RR-7950, INRIA (2012)
4. Bougeret, M., Casanova, H., Robert, Y., Vivien, F., Zaidouni, D.: Using group replication for resilience on exascale systems. Tech. Rep. 265, LAWNS (2012)
5. Bouteiller, A., Bosilca, G., Dongarra, J.: Redesigning the message logging model for high performance. *CCPE* 22(16), 2196–2211 (2010)
6. Buntinas, D., Coti, C., Herault, T., Lemarinier, P., Pilard, L., Rezmerita, A., Rodriguez, E., Cappello, F.: Blocking vs. non-blocking coordinated checkpointing for large-scale fault tolerant MPI protocols. *FGCS* 24(1), 73 – 84 (2008)
7. Cappello, F., Geist, A., Gropp, B., Kalé, L.V., Kramer, B., Snir, M.: Toward exascale resilience. *IJHPCA* 23(4), 374–388 (2009)
8. Davies, T., Karlsson, C., Liu, H., Ding, C., , Chen, Z.: High Performance Linpack Benchmark: A Fault Tolerant Implementation without Checkpointing. In: 25th ICS. pp. 162–171. ACM (2011)
9. Dongarra, J., Beckman, P., et al.: The international exascale software roadmap. *IJHPCA* 25(11), 3–60 (2011)
10. Du, P., Bouteiller, A., et al.: Algorithm-based fault tolerance for dense matrix factorizations. In: 17th SIGPLAN PPOPP. pp. 225–234. ACM (2012)
11. Fagg, G., Dongarra, J.: FT-MPI: Fault tolerant MPI, supporting dynamic applications in a dynamic world. In: 7th EuroPVM/MPI. LNCS, vol. 1908, pp. 346–353. Springer (2000)
12. Gabriel, E., et al.: Open MPI: Goals, concept, and design of a next generation MPI implementation. In: 11th EuroPVM/MPI. LNCS, vol. 3241, pp. 353–377. Springer (2004)
13. Gropp, W., Lusk, E.: Fault tolerance in message passing interface programs. *IJHPCA* 18, 363–372 (2004)
14. Hadzilacos, V., Toueg, S.: *Distributed systems* (2nd ed.). chap. Fault-tolerant broadcasts and related problems, pp. 97–145. ACM/Addison-Wesley (1993)
15. Huang, K., Abraham, J.: Algorithm-based fault tolerance for matrix operations. *IEEE Transactions on Computers* 100(6), 518–528 (1984)
16. Hursey, J., Graham, R.L., Bronevetsky, G., Buntinas, D., Pritchard, H., Solt, D.G.: Run-through stabilization: An MPI proposal for process fault tolerance. In: 18th EuroMPI. LNCS, vol. 6690, pp. 329–332. Springer (2011)
17. Hursey, J., Naughton, T., Vallee, G., Graham, R.L.: A log-scaling fault tolerant agreement algorithm for a fault tolerant MPI. In: 18th EuroMPI. LNCS, vol. 6690, pp. 255–263. Springer (2011)
18. Lusk, E., Chan, A.: Early experiments with the OpenMP/MPI hybrid programming model. In: 4th IWOMP, LNCS, vol. 5004, pp. 36–47. Springer (2008)
19. Mohan, C., Lindsay, B.: Efficient commit protocols for the tree of processes model of distributed transactions. In: SIGOPS OSR. vol. 19, pp. 40–52. ACM (1985)
20. Sterling, T.: HPC in phase change: Towards a new execution model. In: HPCCS – VECPAR 2010, LNCS, vol. 6449, pp. 31–31. Springer (2011)