

Detection and Correction of Silent Data Corruption for Large-Scale High-Performance Computing

David Fiala, Frank Mueller
North Carolina State University
Raleigh, NC
{dfiala, fmueller}@ncsu.edu

Christian Engelmann
Oak Ridge Natl Lab
Oak Ridge, TN
engelmannc@ornl.gov

Rolf Riesen
IBM Ireland
Dublin, Ireland
rolf.riesen@ie.ibm.com

Kurt Ferreira, Ron Brightwell
Sandia Natl Labs
Albuquerque, NM
{kbferre, rbrigh}@sandia.gov

Abstract—Faults have become the norm rather than the exception for high-end computing clusters. Exacerbating this situation, some of these faults remain undetected, manifesting themselves as silent errors that allow applications to compute incorrect results.

This paper studies the potential for redundancy to detect and correct soft errors in MPI message-passing applications while investigating the challenges inherent to detecting soft errors within MPI applications by providing transparent MPI redundancy. By assuming a model wherein corruption in application data manifests itself by producing differing MPI messages between replicas, we study the best suited protocols for detecting and correcting corrupted MPI messages.

Using our fault injector, we observe that even a single error can have profound effects on applications by causing a cascading pattern of corruption which in most cases spreads to all other processes. Results indicate that our consistency protocols can successfully protect applications experiencing even high rates of silent data corruption.

I. INTRODUCTION

In High-End Computing (HEC), faults have become the norm rather than the exception for parallel computation on clusters with 10s/100s of thousands of cores. Past reports attribute the causes to hardware (I/O, memory, processor, power supply, switch failure etc.) as well as software (operating system, runtime, unscheduled maintenance interruption). In fact, recent work indicates that (i) servers tend to crash twice a year (2-4% failure rate) [1], (ii) 1-5% of disk drives die per year [2], (iii) DRAM errors occur in 2% of all DIMMs per year [1], which is more frequent than commonly believed, and (iv) large scale studies indicate that simple ECC mechanisms alone are not capable of correcting a significant number of DRAM errors [3].

In response, long-running applications on HEC installations are required to support the checkpoint/restart (C/R) paradigm to react to faults. This is particularly critical for large-scale jobs; as the core count increases, so does the overhead for C/R, and it does so at an exponential rate. This does not come as a surprise as any single component failure suffices to interrupt a job. As we add system components (multicore chips, memory, disks), the probability of failure combinatorially explodes.

Prior work [4], [5], [6] has revealed that checkpoint/restart efficiency, *i.e.*, the ratio of useful vs. scheduled machine time, can be as high as 85% and as low as 55% on current-generation HEC systems. Recent work by Sandia [7] shows rapidly decaying useful work for increasing node counts (see Table I).

TABLE I
168-HOUR JOB, 5 YEAR MTBF

# Nodes	work	checkpt	recomp.	restart
100	96%	1%	3%	0%
1,000	92%	7%	1%	0%
10,000	75%	15%	6%	4%
100,000	35%	20%	10%	35%

Only 35% of the work is due to computation for a 168 hour job on 100k nodes with a node MTBF of 5 years while the remainder is spent on checkpointing, restarting and then partial recomputation of the work lost since the last checkpoint. Table II shows that for longer-running jobs or shorter MTBF, useful work becomes *insignificant* as most of the time is spent on restarts.

TABLE II
100K NODE JOB, VARIED MTBF

job work	MTBF	work	checkpt	recomp.	restart
168 hrs.	5 yrs	35%	20%	10%	35%
700 hrs.	5 yrs	38%	18%	9%	43%
5,000 hrs.	1 yr	5%	5%	5%	85%

The most important finding of the Sandia study is that **redundancy in computing can significantly revert this picture**. By doubling up the compute nodes so that every node N has a replica node N', a failure of primary node N no longer stalls progress as the replica node N' can take over its responsibilities. Their prototype, rMPI, provides dual redundancy [7]. And *redundancy scales*: As more nodes are added to the system, the probability for simultaneous failure of a primary N *and* its replica rapidly decreases. Of the above overheads, the recompute and restart overheads can be nearly eliminated (to about 1%) with only the checkpointing overhead remaining — at the cost of having to deploy twice the number of nodes (200,000 nodes in Table I) and four times the number of messages [7]. But once restart and rework overheads exceed 50%, redundancy is actually *cheaper* than traditional C/R at large core counts.

The failure scenarios above only cover a subset of actual faults, namely those due to fail-stop behavior / those detectable by monitoring of hardware and software. Silent data corruption (SDC) is yet a different class of faults, which is

the focus of this work. It materializes as bit flips in storage (both volatile memory and non-volatile disk) or even within processing cores. A single bit flip in memory can be detected (with CRC) and even mitigated with error correction control (ECC). Double bit flips, however, force an instant reboot after detection since ECC cannot correct such faults. While double bit flips were deemed unlikely, the density of DIMMs at Oak Ridge National Lab’s Cray XT5 causes them to occur on a daily basis (at a rate of one per day for 75,000+ DIMMs) [8].

Meanwhile, even single bit flips in the processor core remain undetected as only caches feature ECC while register files or even ALUs typically do not. Significant SDC rates were also reported for BG/L’s unprotected L1 cache [9], which explains why BG/P provides ECC in L1. Nvidia made a similar experience with its shift to ECC in their Fermi Tesla GPUs. Yet, hardware redundancy remains extremely costly [10], [11], [12], [13]

Today, the frequency of bit flips is no longer believed to be dominated by single-event upsets due to radiation from space [14] but is increasingly attributed to fabrication miniaturization and aging of silicon given the increasing likelihood of repeated failures in DRAM after a first failure has been observed [1]. With SDCs occurring at significant rates, not every bit flip results in faults. Flips in stale data or code remain without impact, but those in active data/code may have profound effects and potentially render computational results invalid without ever being detected. This creates a severe problem for today’s science that relies increasingly on large-scale simulations. Redundant computing can detect SDCs where relevant, i.e., when results are impacted. While detection requires dual redundancy, correction is only feasible with triple redundancy. Such high levels of redundancy appear costly, yet may be preferable to flawed scientific results. Triple redundancy is also cheaper than comparing the results of two dual redundant jobs, which would be the alternative at scale given the amount of useful work without redundancy for large systems from Table II. Overall, the state of HEC requires urgent investigation to level the path to exascale computing — or exascale HEC may be doomed as a failure (with very short mean times, ironically).

A. Modeling Redundancy

Elliott *et al.* [15] combine partial redundancy with checkpointing in an analytic model and experimentally validate it. Results indicates that for an exascale-size machine, more jobs can utilize a cluster under redundancy than would be possible with checkpointing without redundancy for the same number of cores, i.e., redundancy increases capacity computing in terms of HPC job throughput.

We used this model in combination with Jaguar’s system MTBF of 52 hours [16] (equivalent to a node MTBF of 50 years) to assess the viability of redundancy. Consider a 128 hour job (without checkpointing). We then assess the time required for such a job without redundancy (1x), dual redundancy (2x) and triple redundancy (3x) at different node counts under weak scaling and with an optimal checkpoint

interval to minimize overall execution (see Fig. 1). At 18,688 nodes (Jaguar), marked as line *C*, single-node (1x) runs are about 7% faster than dual (2x) and 20% faster than triple (3x) redundancy. The problem is that a job at 1x will have no indication if it had been subjected to an SDC. Consider Jaguar’s double bit error rate of once a day again [8], which is silently ignored (to increase system availability) as it cannot be corrected. Scientists will not know if their outputs were affected, i.e., if outputs are flawed (incorrect science problem).

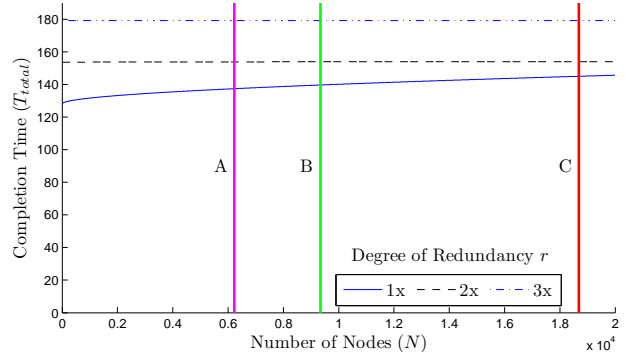


Fig. 1. Modeled Time to Completion with Redundancy

Let us consider dual redundancy at half the node count of Jaguar (line *B* in Fig. 1). In order to ensure absence of SDCs, a user would have to run a single redundant (1x) job twice for a total time of about 280 hours (twice 140) vs. a dual redundant (2x) job at twice the number of nodes (full Jaguar size, line *A*) with 155 hours. Hence, dual redundancy, results in nearly half the wall-clock time if SDC detection is a requirement for verification and validation of a job’s results.

Consider triple redundancy at a third of Jaguar’s node count (line *A*). Running two jobs at 1x takes about 276 hours (twice 138), a dual redundant (2x) job 145 hours and a triple redundant job 180 hours. The output of the two 1x jobs differs, a third run would be required (assuming that two of them produce correct results). If the dual redundant (2x) job detects an error, it also needs to be rerun. The triple redundant (3x) job, in contrast, can correct errors so that no reruns would be needed.

At exascale core counts of one million and a node count of 100,000 (swim lane 1 [17], [18]), dual redundancy would have the lowest cost (lower than single job at 1x). The additional cost of SDC correction at triple redundancy adds another 14% overhead in wall-clock time, with the benefit of no repeated runs for SDCs. This is based on the assumption of Jaguar’s system MTBF (52 hours), even though the MTBF could be much smaller given that double bit errors for a 128 Petabyte system would occur every four minutes (compared to one a day today on Jaguar) [8]. SDC detection, if not correction, may thus become essential at exascale.

B. Contributions

The main contributions of this work are (i) the design of novel silent data correction detection / correction methods

and (ii) a study on the challenges and costs of performing SDC protection using redundancy. By utilizing redundancy, our key to success is to not only rely on reactive resilience requiring restart overheads but to sustain failures with forward computational progress without a need to restart.

Our work makes the following major contributions: (1) We contribute the design and implementation of protocols for SDC detection and correction at the communication layer. (2) We demonstrate the capabilities and assess the cost of redundancy to (a) detect SDC and (b) recover from such corruption in experiments on a real system. While dual redundancy can detect SDCs, triple redundancy can actually correct them through voting. We study the benefits and limitations of the spectrum ranging from no redundancy over dual to triple redundancy in terms of overhead and computing/interconnect resource costs. A key challenge is to limit the overhead for SDC detection by reducing the relevant footprint of computational results, which we explore. (3) We assess the resilience of HEC jobs to faults through injection. Hardware and software failures are studied through injection on an actual cluster. (4) We develop a live SDC tracking and reporting framework to investigate the effects of SDCs on applications in terms of their rate of taint (corruption) progression spreading from node to node via MPI communication. Further, we use this framework to evaluate several application responses to fault injection and classify three types of observed behavior that result in invalid data being generated. In summary, this work contributes to fault detection and recovery by significantly advancing existing techniques by controlling levels of redundancy intervals in the presence of hardware and software faults.

II. DESIGN

This work presents RedMPI, an MPI library that is capable of both detecting and correcting SDC faults. RedMPI creates “replica” MPI tasks for each “primary” task and performs online MPI message verification intrinsic to existing MPI communication. The replicas compare received messages, or hashes, from multiple senders and can thus detect if a process’s communication data has been corrupted.

RedMPI can run in double redundant mode and detect divergent messages between replicas. Such messages are indicative of corruption due to the fact that replicas will be run in a deterministic manner. When RedMPI is run in triple redundant mode, it gains the additional potential to also correct faulty messages from a corrupted replica. RedMPI supports additional levels of redundancy for environments where multiple near-simultaneous faults can occur during data transmission. A voting algorithm is used to determine which of the received messages are correct and should be used by all receivers.

To detect SDCs, RedMPI solely analyzes the content of MPI messages to determine divergence between replicas during communication. Upon divergence, the result deemed to be invalid is discarded on the receiver side and transparently replaced with a known “good” value from another replica.

A different SDC detection approach would be to constantly compare the memory space of replicas’ processes and compare

results. Such an approach suffers from excessive overhead due to constant traversals of large memory chunks, overhead due to global synchronization to ensure that each process is paused at the exact same spot during a memory scan, and the communication required for replicas to compare their copy of each memory scan while looking for differences. In this case, if corruption is detected, it is not feasible to correct the memory while the application is running as this could interfere with application-side writes to the same memory region. This, in turn, could necessitate a rollback of all tasks to the last “good” checkpoint (assuming that checkpointing was also enabled).

By instead focusing on the MPI messages themselves, we have cut our search area down to only data that is most critical for correctness of an MPI application; i.e., we argue communication correctness is a necessary (but not sufficient) condition for output correctness. Moreover, should an SDC occur in memory that is not immediately communicated over MPI, the fault is eventually detected as the corrupted memory may later be accessed, operated on, and finally transmitted. The same principle holds true for data that became corrupted while residing in a buffer or any other place in memory. If the SDC is determined to eventually alter messages, then RedMPI detects it during transmission, independent of when, where or how the SDC originated.

It is very important to note that RedMPI is designed to protect an entire application from SDC by using replication. RedMPI is not designed to protect an interconnect. By assuming that an application’s most critical data is communicated during/after computation, we have effectively reduced the scope to data that gets communicated and may be compared between replicas to ensure consistency. A process receiving corrupted data that affects important calculations will eventually result in message correction so that uncorrupted replicas are guaranteed to have received correct data, only.

A. Point-to-Point Message Verification

The core of RedMPI’s error detection capabilities are designed around a reliable, verifiable point-to-point communication protocol. Specifically, a point-to-point message (e.g., `MPI_Isend`) sent from an MPI process must be identical to the message sent by other replicas for any given rank. Upon successful receipt of a message, the MPI application is assured that the message is valid (not corrupted).

Internally, a verification message may take the form of a complete message duplicate that is compared byte-by-byte. Alternatively, since MPI messages may be large, it is in many cases more efficient to create a unique hash of the message data and use the hash itself for message verification to reduce network bandwidth. Message data verification can be performed at either the sender or the receiver.

Let us first consider the case of sender-side verification. To perform verification at the sender, all of the replicas need to send a message to communicate with each other and verify their content (through some means) before sending the verified data to the receiving replicas. However, this approach incurs

added latency and overhead for each message sent due to the time taken to transmit between replicas and to internally verify messages. Additionally, it is best to optimize for the critical path, i.e., for the case that a sent message is not corrupted and that all senders have matching data. A sender-side approach is subject to additional overhead for every message sent at both sending and receiving nodes. Specifically, while every sent message is treated as suspect, the time required for the senders to agree that each of their own buffered messages is correct represents the time lost on the receiver side before the application can proceed. For this reason, RedMPI's protocols use receiver-side verification resulting in faster message delivery with considerably reduced message latency.

B. Assumptions

For brevity, several aspects of the RedMPI implementation are not included here, but are published in [19]. Many technical contributions and challenges are available within that document and include RedMPI's sphere of protection, how collectives/wildcards/non-deterministic MPI operations are entirely supported, and details of how we provide transparent MPI redundancy for nearly all MPI-1 functions without application modifications. RedMPI does not protect MPI I/O functionality, but orthogonal work [20] could be combined with RedMPI to also cover this aspect.

C. Implementation Notes

RedMPI provides the capability of soft error detection for MPI applications by online comparison of results of identical replica MPI processes. To an MPI developer, the execution of replica processes of their original code is transparent as it is handled through MPI introspection within the RedMPI library. This introspection is realized through the *MPI profiling layer* that intercepts MPI function calls and directs them to RedMPI. The profiling layer provides a standard API allowing libraries to wrap MPI calls and add additional or replacement logic for API calls.

To understand how RedMPI functions internally, it is first important to understand how redundancy is achieved within RedMPI. When launching an MPI job with RedMPI, some *multiple* of the original number of desired processes needs to be launched. For example, to launch an MPI job that normally requires 128 processes it would instead require 256 or 384 processes for dual or triple redundancy, respectively. RedMPI handles redundancy internally and provides an environment to the application that appears to only have the originally required 128 processes. (i.e., `Comm_rank`, `Comm_size`, `Send`, `Recv`, etc., all appear to operate normally to the running MPI application even with redundancy.)

1) *Message Corruption Detection & Correction*: RedMPI's first receiver-side protocol, All-to-all, supports both message verification and message voting to ensure that the receiver discards corrupted messages. The All-to-all method requires that each MPI message sent is transmitted from all sender replicas to each and every receiver replica. Thus, for triple redundancy, each sender sends three messages where one

message goes to each replica receiver.

Following a message transmission or reception, an MPI application usually completes these requests with an `MPI_Test` or `MPI_Wait`. RedMPI interposes these functions as it needs to test not just the single `MPI_Request`, but rather impose a test for multiple requests corresponding to sends/receives from all replicas. The `MPI_Request` is looked up and the test or wait is performed on all outstanding requests. If the test or wait was performed on a request from an `MPI_Isend`, then no further action from RedMPI is required once the requests complete. Only a request from an `MPI_Irecv` requires extra steps to verify matching message reception from each replica.

For dual redundancy, messages received from replicas should match. If they do not, then RedMPI reports the corruption to the user and may terminate the application. Triple redundancy or better uses voting to determine which sending replica has transmitted unmatching, corrupt data. The MPI application's receive buffer will always receive a copy of a correct message per the results of the voting.

RedMPI's primary mode of operation is its second, more efficient protocol: `MsgPlusHash` (message plus hash). The `MsgPlusHash` corruption detection and correction method provides a key performance enhancement over the All-to-all method by vastly reducing the total data transfer overhead per message and the number of messages in the general case. Similar to the All-to-all method, `MsgPlusHash` performs message verification solely on the receiver end. The critical difference is that `MsgPlusHash` sends one copy of a message originating from an MPI transmission in addition to a very small hash message. This change in protocol allows each sending replica to transmit their message only once, while the additional hash message is later used to verify each receiver's message. `MsgPlusHash`'s contribution is a reduction of messages and thus bandwidth required from n^r to simply $n * r$ (plus a hash) where n is the number of messages sent and r is the degree of replication.

To detect message corruption, the minimum requirement is a comparison between two different sources. Additionally, the most likely scenario (critical path) is for corruption to not exist. The `MsgPlusHash` method takes full advantage of these facts by only receiving a single copy of any message transmitted plus a hash from an alternate replica. From an efficiency standpoint, it is not necessary to send two full messages since a hash suffices to verify data correctness without imposing overheads of full message retransmission. Once the full message is received, a hash of the message is generated at the receiver and compared with a hash from a different replica. In the likely event that hashes match, the receiver can be assured that its message is correct, i.e., no corrective action is taken.

As shown in Figure 2, each sender replica must calculate where to send its message and where to send a hash of its message. The actual message's destination is simply calculated by finding the receiver with the same replica rank as the sender. The hash message's destination is calculated by taking the sender's replica rank and adding one (supporting wrap-around

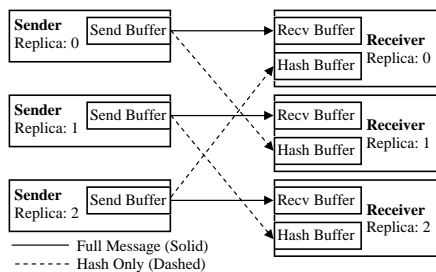


Fig. 2. MsgPlusHash Method

to 0).

RedMPI employs an advanced discovery and message correction protocol to identify and fix corrupt message data in the event that a replica receives a message and hash that do not match. To conserve space, the reader is referred to [19] for details.

III. FAULT INJECTOR DESIGN

As the research goals of this work include detecting and protecting applications from silent data corruption, an integrated fault injection tool is required to evaluate the effectiveness of RedMPI to detect and correct memory errors during execution. Additionally, the same fault injector can later be used to monitor the adverse effects of SDC on running applications when RedMPI is not actively protecting them, i.e., we specifically injector errors and allow them to propagate.

To experimentally determine the effect of corruption and verify corrective actions, our fault injector was designed to produce data corruption in a manner resembling naturally occurring faults. Namely, single bit flips undetected by ECC are of interest (e.g., within an arithmetic-logic unit of a processor) when their effects eventually propagate into a message transmission over MPI. Alternatively, also of interest is SDC due to multiple bit flips in main memory resulting in a corrupted bit pattern that ECC is unable to detect / correct.

Our fault injector, which is built to co-exist with RedMPI, specifically targets MPI message send buffers to ensure that each injection actually impacts the MPI application while simultaneously reaching message recipients. When activated, the fault injector is given a frequency of $1/x$ during launch, which is the probability that any single message may become corrupted. By using a random number generator with a state internal to RedMPI (without effect on the MPI application), the injector randomly picks messages to corrupt. Once targeted for corruption, RedMPI selects a random bit within the message and flips it prior to sending it out. RedMPI is agnostic to the data type of the message, i.e., the injector calculates the total number of bits within the entire message regardless of type or count before picking a bit to flip.

Note that not only does the fault injector flip a bit in the send buffer, but it actually modifies the application’s memory directly. If the MPI application accesses the same memory again, further calculations based on that data will be invalid

with a high probability of causing further divergence from non-corrupted replicas.

A. Targeted Fault Injections

Memory corruption faults may also be specifically targeted to occur within specific sets of replicas, MPI ranks, application timesteps, or frequency. For instance, as we will later investigate the effects of SDC on our experimental applications when SDC errors are allowed to propagate without any protection enabled, our targeted faults will be limited to only one set of replicas such that in dual redundancy half of the nodes may serve as “control” replicas that never experience faults while the other “experiment” replicas will receive faults. By modifying MPI applications to report back to RedMPI whenever they reach a new timestep, we can also target faults to occur at very specific points in execution such as defining a desired timestep for an SDC injection to occur.

IV. EXPERIMENTAL FRAMEWORK

We deployed RedMPI on a medium sized cluster and utilized up to 96 nodes for benchmarking and testing. Each compute node consists of a 2-way SMPs with AMD Opteron 6128 (Magny-Cours) processors of 8 cores per socket (16 cores per node) with 32 GB RAM per node. Nodes are connected via 1Gbps Ethernet for user interactions and management. MPI transport is provided by a 40Gb/s InfiniBand fat tree interconnect. To maximize the compute capacity of each node, we ran up to 16 processes per node.

When launching RedMPI jobs, we map replica processes so that they do not reside on the same physical nodes. This type of mapping is preferred as a fault on a node will not affect multiple replicas of the same process simultaneously (i.e., due to localized power failures for a whole rack).

A. Time Overhead Experiments

RedMPI allows applications to utilize transparent 2x or 3x redundancy. While the physical cost of redundancy is known (2x or 3x the number of tasks), the additional cost in terms of wall-clock time should be investigated to determine what types of costs are expected, if any. To determine the cost of redundancy in terms of time, we run a variety of applications demonstrating differing scaling, processor counts, communication patterns, and problem sizes.

For our timing experiments, we solely report benchmark results for the MsgPlusHash SDC method as it provides a more efficient communication protocol than All-to-all by design. To provide meaningful metrics, each experiment assesses the run time for regular, unaltered Open MPI (referred to as $1x$ in the results tables), RedMPI with dual redundancy (2x), and RedMPI with triple redundancy (3x). Note that the size reported is the size without redundancy in all experiments and results. Hence, a size 512 job in our results with triple redundancy is actually running across 1536 processors.

We assess both strong and weak scaling when evaluating overheads associated with RedMPI. For each weak scaling application, the input data size remains constant for each process no matter how many processes are run. In contrast,

TABLE III
LAMMPS INPUT CHAIN.SCALED

Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV
128	240.5	241.34	242.54	-3.8%	-3.3%
256	244.39	244.61	245.25	0.1%	0.4%
512	250.93	251.89	256.11	0.4%	2.1%

TABLE VI
HPCCG (USES MPI WILDCARDS)

Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV
128	99.79	99.76	125.75	0.0%	26.0%
256	99.64	128.83	131.02	29.3%	31.5%
512	126.36	146.19	152.26	15.7%	20.5%

TABLE IX
NPB FT

Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV
32-C	117.45	117.95	118.68	0.43%	1.05%
64-C	68.82	68.62	71.77	-0.29%	4.29%
128-D	222.75	228.76	234.97	2.70%	5.49%

TABLE IV
LAMMPS INPUT CHUTE.SCALED

Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV
128	137.50	138.38	139.01	0.6%	1.1%
256	138.26	140.43	140.00	1.6%	1.3%
512	139.19	140.22	140.67	0.7%	1.1%

TABLE VII
NPB CG

Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV
128-D	201.42	205.87	215.51	2.2%	7.0%
256-D	127.21	132.61	136.64	4.2%	7.4%
512-D	70.10	77.54	83.67	10.6%	19.4%

TABLE X
NPB LU

Size	1x [sec]	2x [sec]	3x [sec]	2x OV	3x OV
128-D	361.78	379.90	375.44	5.0%	3.8%
256-D	179.78	191.97	195.55	6.8%	8.8%
512-D	102.90	115.07	121.01	11.8%	17.6%

TABLE V
SWEEP3D

Size	1x [s]	2x [s]	3x [s]	2x OV	3x OV
128	390.30	389.49	393.05	-0.2%	0.7%
256	428.17	427.53	431.20	-0.1%	0.7%
512	488.08	488.93	494.09	0.2%	1.2%

TABLE VIII
NPB EP

Size	1x [s]	2x [s]	3x [s]	2x OV	3x OV
128-D	72.31	72.63	72.74	0.4%	0.6%
256-E	579.94	581.02	581.27	0.2%	0.2%
512-E	289.80	290.83	291.30	0.4%	0.5%

TABLE XI
NPB MG

Size	1x [s]	2x [s]	3x [s]	2x OV	3x OV
128-E	339.17	340.41	429.67	0.4%	26.7%
256-E	168.56	170.68	171.48	1.3%	1.7%
512-E	66.97	68.35	69.29	2.1%	3.5%

strong scaling applications have a constant problem size for a given class that varies the amount of input data each process receives when jobs of differing sizes are run. In effect, we expect strong scaling applications to reduce the amount of data and computation required per process as the number of processors increases.

Our test suite of weak scaling applications includes LAMMPS, ASCII Sweep3D, and HPCCG. LAMMPS is a popular molecular dynamics code that we evaluate with two different problems, “chain” and “chute”. Sweep3D is a neutron transport code. Finally, HPCCG is a finite elements application from the Sandia National Labs Mantevo Project. It was chosen because of its use of `MPI_ANY_SOURCE` to demonstrate RedMPI’s capability to handle non-deterministic MPI operations.

The strong scaling NAS Parallel Benchmarks (NPB) are also evaluated with varying problem class sizes and number of processes. We use the NPB suite to demonstrate how varying the communication-to-computation ratio affects RedMPI in some cases.

V. RESULTS

Tables III-XI report execution time for the benchmarked applications. Every application was run with three different MPI sizes. For all of the cases except one, we conducted experiments with 128, 256, and 512 processors for the baseline. The uninstrumented (no RedMPI) version of each application is shown under the *1x* column while the *2x* and *3x* columns represent dual and triple redundancy, respectively, under RedMPI. The final two columns represent the percent overhead incurred by adding dual or triple redundancy relative to the baseline. Runs with redundancy use two or three times as many processes as the uninstrumented baseline runs. Performance is subject to cache effects when running the same application with RedMPI. This effect may vary between degrees of redundancy. This is evident for results that indicate small negative overheads (speedup under redundancy) when averaged, such as in Table III.

We first analyze the runtime results of the weak scaling applications in Tables III-VI. LAMMPS with input *chain* was run with a dataset size of $32 \times 40 \times 20$ for 512 processors

and scaled down proportionally for 256 and 128 processors, which explains the relatively consistent runtimes. Likewise, LAMMPS with input *chute* was given a dataset size of 320×88 for 512 processors and also scaled proportionally. Sweep3D had an input size of $320 \times 40 \times 200$ and HPCCG had an input size of $400 \times 100 \times 100$.

These applications performed very well with RedMPI, i.e., in most cases the RedMPI overhead was not perceptible due to a well balanced communication-to-computation ratio that weak scaling allowed us to retain despite increasing the number of processors as we scaled up the benchmarks with RedMPI.

To demonstrate the effectiveness of RedMPI’s wildcard support, HPCCG was chosen as it makes use of `MPI_ANY_SOURCE` receives. Since RedMPI requires special handling for wildcards, the overheads incurred may vary based on how long it takes the replicas to receive an envelope message that resolves the wildcard. When wildcard resolution is completed quickly, very little performance penalty is seen as in the *2x* results with size 128 (see Table VI). Conversely, when wildcard resolution takes a relatively long time, then RedMPI forces MPI to receive all messages in an unexpected queue. Delaying message reception can potentially degrade performance, and MPI wildcards are recommended to be avoided when possible.

The remaining NPB benchmarks are strong scaling applications subjected to input classes C, D, and E, where E is the largest size. For each class, data is equally distributed across all MPI processes. For example, the CG and LU benchmarks were both run with the same class size *D* for 128, 256, and 512 size jobs. We can see that as the number of processes increases, the baseline runtime decreases since there is less computation per process to perform. In turn, as the computation is decreased per process the amount of communication incurred does, in fact, increase. Tables VII and X clearly demonstrate how the overhead of RedMPI increases as the per-process communication overshadows the per-process computation. Hence, to keep RedMPI overheads reasonable, it is important to choose input classes such that the ratio of communication-to-computation is balanced, e.g.,

as seen for input sizes for EP and FT. For FT, we were unable to run class E with size 256 and 512 because our experimental setup did not have enough memory available to hold class E. Thus, we chose to run FT with smaller class sizes and report smaller runs.

Overall, RedMPI's runtime overheads are modest for well-behaved applications that can be scaled to exhibit a fair communication-to-computation ratio.

VI. FAULT INJECTION STUDIES

RedMPI's fault injector provides two key opportunities for specifically analyzing silent data corruption faults within the scope of running MPI applications. We will use RedMPI to answer these questions:

(1) **Propagation:** Does SDC affect applications messages and correctness when no protection mechanisms (such as RedMPI's voting) are available? How quickly do SDC injections propagate to other processes via communication? Do corrupted processes further disrupt other processes in a cascading manner by sending invalid, divergent MPI messages as compared to the correct execution of a job?

(2) **Protection:** When utilizing triple redundancy with RedMPI, are SDCs successfully detected and corrected? Do applications still complete with correct answers even in the face of SDC injections?

A. SDC Propagation Study

Our first study investigates whether leaving message data unprotected in the face of SDCs does in fact lead to incorrect results. This happens when a single SDC injection in one process will later spread to other nodes causing an overwhelming cascade of invalid data.

As described in Section III-A, we run RedMPI with dual redundancy and assume two sets of replicas. The first set of replicas is a control set and will not receive any SDC injections. The control set will execute normally and should produce correct results upon completion. Our second set of replicas is the test set, which becomes the victim when SDC faults are injected. *Further, during these experiments we also purposefully disable RedMPI's corrective capabilities.* RedMPI still detects divergent messages between the control and test set replicas, but allows application progress to continue. RedMPI tracks live statistics on applications running in this environment such as:

- which processes receive SDC injections;
- which processes send bad messages and where;
- which processes receive corrupt messages and how that corruption further spreads to and from nodes that were indirectly tainted by a bad message; and
- aggregate data at the granularity of a single application-defined timestep.

This type of reporting is made possible by redundancy and is considered a new technique for live application analysis and correctness since we do not need to actually log data for later viewing; this is advantageous for long-running applications or when bandwidth is high and logging to disk for offline analysis is undesirable.

Figures 4-10 show how SDC injection(s) spread causing bad messages and cascading tainting of other processes via reception of corrupt messages. The x-axis of these graphs denotes progress in the form of application-reported timesteps. The blue line correlates with the left vertical axis, which denotes the number of bad messages that were cumulatively received by all MPI processes in any given timestep. The gray filled areas match to the right axis and denote the number of MPI processes that receive an SDC injection themselves (direct tainting) or become corrupt indirectly (indirect tainting/light gray) by receiving a corrupt message from a previously corrupted sender themselves. Combined, both of the gray filled areas indicate how many of the MPI processes operating in a single job have become corrupt over time.

For all of the experiments except NPB FT (Figure 10), only a single SDC injection was performed. NPB FT was subjected to 4 SDC injections. All SDC injections were randomly inserted using the aforementioned RedMPI fault injector. In all cases we ran the experiments with 64 MPI processes under dual redundancy (64 control set processes + 64 test set processes).

These 7 graphs indicate three disjoint trends in response to SDC injections. First, the *progressive* trend is characterized by Figures 4 and 5 as a single injection that does not spread immediately. Instead, *progressive* applications often communicate with their grid neighbors resulting in a sphere of corruption that grows outward. For example, Figure 3 shows a heatmap indicating nodes (both on the x-and y-axes) that communicate tainted messages with one another due to an SDC. The figures depict timesteps 4 (one timestep after first injection), 58 (more indirectly tainted messages), 112 (most nodes tainted), and 221 (all nodes tainted) after the SDC injection. As the application progresses, the final heatmap shown is actually a heatmap of normal communication. This indicates that eventually virtually all communication has succumbed to invalid data due to a single SDC injection.

The next trend we identify is the *immediate/explosion* case. Of the experiments reported, Figures 6-9 fall in this category. For these applications we noticed that a single injection resulted in corruption that spread across all nodes usually within two timesteps of the initial SDC. The most common reason for the *immediate* trend is the use of collectives or some communication pattern that tends to exchange messages between all nodes in a short period of time. Heatmaps of this trend are not provided as they essentially mirror the full communication pattern of the application almost immediately after the injection occurs.

Our third identified trend is the *localized* case. In Figure 10, NPB FT is one such application where injections result in just a few invalid messages. This occurs when the corrupted data is neither reused nor retransmitted, which in turn keeps the sphere of corruption relatively isolated to the processes that were tainted by a direct SDC injection. In this experiment we targeted FT with multiple injections to demonstrate how applications that fit the *localized* trend typically do not result in large aggregates of tainted messages or nodes. Although we

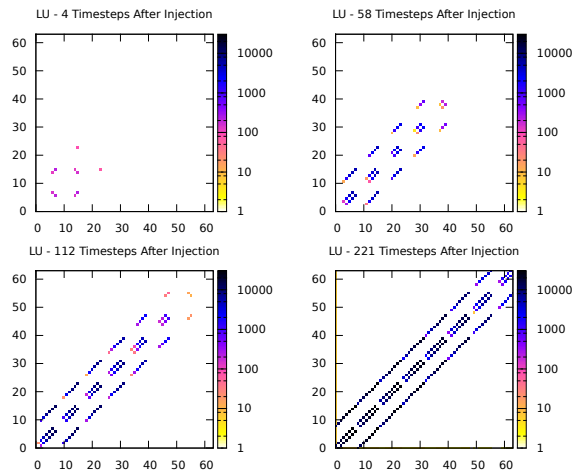


Fig. 3. NPB LU Corrupted Communication Patterns

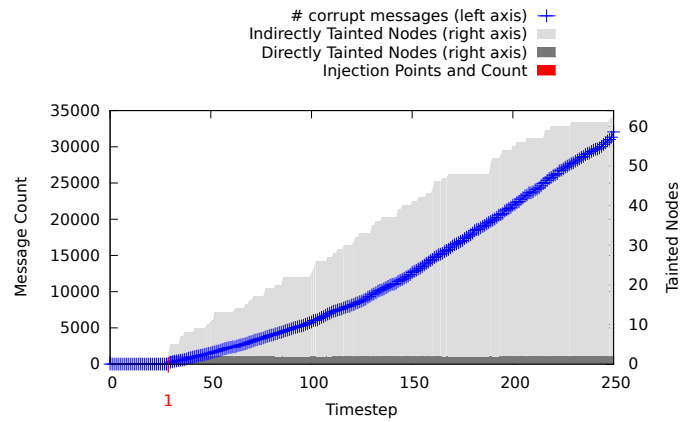


Fig. 4. NPB LU Overview of Corrupt Nodes and Messages

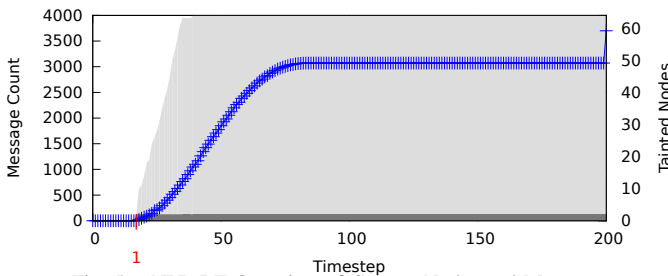


Fig. 5. NPB BT Overview of Corrupt Nodes and Messages

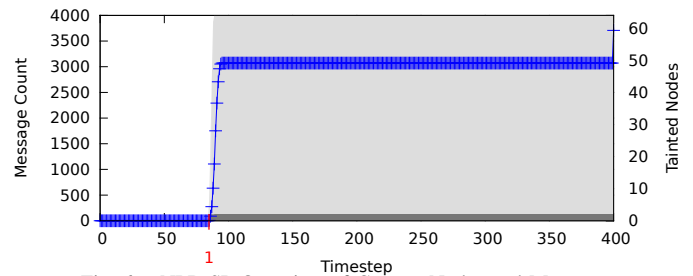


Fig. 6. NPB SP Overview of Corrupt Nodes and Messages

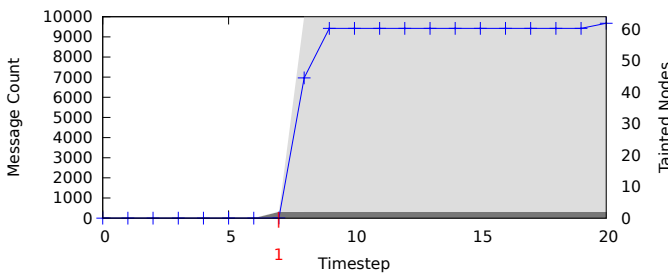


Fig. 7. NPB MG Overview of Corrupt Nodes and Messages

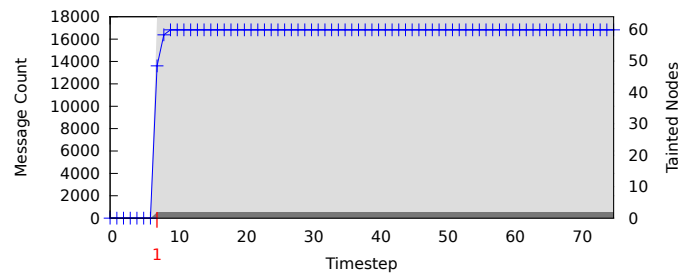


Fig. 8. NPB CG Overview of Corrupt Nodes and Messages

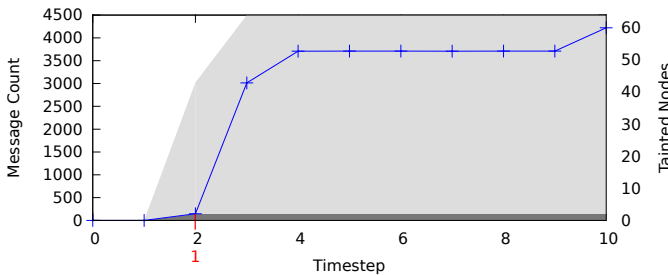


Fig. 9. ASCI Sweep3D Overview of Corrupt Nodes and Messages

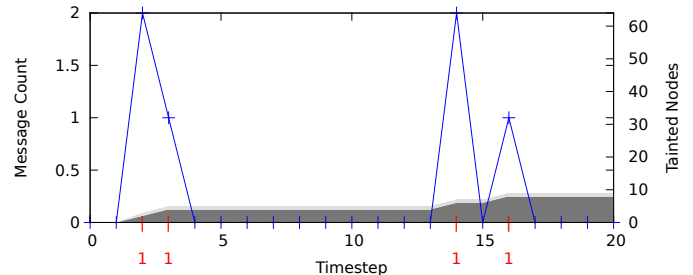


Fig. 10. NPB FT Overview of Corrupt Nodes and Messages

did not receive a high degree of corrupt messages or spreading for FT, it is important to note that all benchmarks (including NPB FT) failed to pass their internal verification or complete with correct results that matched their “control” counterparts.

In summary, by observing that just a single SDC injection can induce a profound effect on all communicating processes, we conclude that protecting applications at the MPI message level is an appropriate method to detect, isolate, and prevent further corruption. Had RedMPI’s protection not been purposefully disabled for this study, then all of the SDC injections would have been isolated from spreading and no bad messages could have been received by other processes.

B. RedMPI SDC Protection Study

To gauge SDC sustainability when RedMPI is active with redundancy, we inject faults into the running benchmarks to determine if the faults are detected and if correction succeeds. Additionally, we experimentally determine if the fault corrections allow the benchmarks to complete their self-verification process successfully.

Next, we analyze the effectiveness of SDC detection and correction protocols. We ran the fault injector with a corruption frequency of $1/5,000,000$ messages to ensure a relatively high likelihood for an injection while running the CG benchmark

with 64 processes (virtual ranks) and a replication degree of three (192 physical processes). Note that we restricted the intentional corruption injections to only occur on replicas with a replica rank of zero to control the experiment. This ensures that at least two of the three replicas do not receive an injection in all cases. Thus, voting results in a valid outcome so that invalid messages can be corrected. During ten experiments with a frequency of $1/5,000,000$, we encountered one occasion with two injections, four occasions with a single injection, and five occasions without injections. In every run except one, the corruption resulted in a single bad message that was successfully detected and corrected by the receiving replicas. In one event, a single injection cascaded resulting in 6,242 bad messages originating from the corrupted sender. Nevertheless, the receiving replicas were able to correct the messages as they arrived. Eventually, the corrupted node ceased to send corrupted messages as the application finished traversing data structures until the fault was no longer touched. In these experiments, the applications progressed until completion and successfully passed their built-in verification at the end of processing.

Following that experiment, we performed injections with a frequency of $1/2,500,000$ in another ten runs that were not limited to a single group of replicas. By doubling the odds for an injection and removing the process selection restriction, we detected a significantly larger number of faults. On average, we received 2.5 injections and several thousand invalid messages per run as a result. Nonetheless, RedMPI carried all but two runs to a successful completion with verification. Of the two runs that failed, we observed that when two of the three replicas simultaneously transmitted a corrupt message over RedMPI, it was detected. The voting process is then forced to fail. In this case, RedMPI aborts because voting becomes impossible with three unique messages and hashes. Statistically, as more processes are added to a job, the likelihood of two processes with the same virtual rank becoming corrupt decreases. Therefore, observing results similar to this particular controlled experiment decreases as job sizes and replication degrees increase. Nonetheless, the longer a process remains in an invalid state (i.e., sending bad messages), the longer the correction features of RedMPI are impaired. However, it is important to note that RedMPI still forces a job to abort if it does detect corruption across two or more nodes while utilizing triple redundancy.

Performance of SDC correction has proved to be quite efficient. During SDC correction overhead experiments, we discovered that with as few as three injections we were able to produce nearly 100,000 invalid messages from corrupted senders. The receiving replicas were able to successfully detect and correct each invalid message while effectively generating no perceived overhead. In fact, while running 20 experimental iterations to gauge the protocol overhead of correcting MsgPlusHash messages during injection, the runtime was 0.31 seconds less than the original experiment runtime without fault injection on average.

Realistically, we do not expect to encounter such a high

number of naturally occurring SDCs for a small environment such as our benchmarking cluster. The actual overhead incurred due to SDC correction is a function of the number of invalid messages, and the number of invalid messages is highly dependent on the data reuse patterns of an MPI application.

As empirical evidence of the success of RedMPI, we discovered that RedMPI was detecting and correcting MPI transmission errors during our timing benchmarks even though we had not activated our fault injection module. While investigating, we learned that RedMPI had been properly detecting and correcting faulty memory that was later confirmed to be producing errors, which could not be corrected by ECC alone. These problems were occasionally (but not always) visible through the Linux EDAC monitoring module of the memory controller. Interestingly, EDAC was unable to consistently detect all of the errors that RedMPI detected. Using RedMPI, we discovered which of the 1536 compute cores in the experiment was faulting and were able to reproduce similar experiments that consistently produced faulty MPI messages on this node before removing it from the production system. Without RedMPI, this failing hardware may have gone unnoticed for some time.

VII. RELATED WORK

Since the early 1990s [21], fault tolerance in large-scale HPC systems is primarily assured through application-level checkpoint/restart (C/R) to/from a parallel file system. Support for C/R at the system software layer exists, e.g., via Berkeley Lab Checkpoint Restart (BLCR) [22] or diskless C/R via Scalable C/R (SCR) [23]. Message logging, algorithm-based fault tolerance, proactive fault tolerance, and Byzantine fault tolerance have all been researched in the past. Redundancy in HPC, as showcased in this paper, has only been recently explored.

Historically, the primary defense against SDC has been ECC in memory. In today's memory modules and processors, single-error correction (SEC) double-error detection (DED) ECC protects against single event upset (SEU) bit flips as well as single event multiple upset scenarios. Chipkill offers additional protection against wear-out and complete failure of a memory module chip by spanning ECC across chips but Bose Chaudhuri-Hocquenghem (BCH) encoding provides better energy-delay characteristics [24]. Software redundancy may provide more extensive SDC protection, especially considering the expected increase in SECDED ECC double-error rates.

Pure software-based solutions [25] try to protect against memory corruption without extending hardware ECC. However, they cannot provide perfect coverage to all memory and are subject to job failure if just a single process terminates due to a fault. In contrast, redundancy for SDC correction survives single process faults more gracefully.

Studies primarily done at Los Alamos National Laboratory focused on analyzing the probability and impact of silent data corruption in HPC environments. One investigation [26] showed that a Cray XD1 system with an equivalent number of processors as the ASCI Q system would experience one

SDC event every 1.5 hours. Another study [27] at Lawrence Livermore National Laboratory investigated the behavior of iterative linear algebra methods when confronted with SDC in their data structures. Results show that linear algebra solvers may take longer to converge, not converge at all, or converge to a wrong result.

Modular redundancy has been used in information technology, aerospace and command & control systems [28]. Recent software-only approaches [29], [30] focused on thread-level, process-level and state-machine replication to eliminate the need for expensive hardware. The sphere of replication [31] concept describes the logical boundary of redundancy for a replicated system. Components within such a sphere are protected; those outside are not.

Recent work [32] studied the impact of deploying redundancy in HPC systems. Redundancy can significantly increase system availability and correspondingly lower the needed component reliability. Redundancy applied to a single computer decreases the MTTF of each replica by a factor of 100-1,000 for dual redundancy and by 1,000-10,000 for triple redundancy without lowering overall system MTTF. If a failed replica is recovered through rebooting or replacing with a hot spare, replica node MTTF can be lowered by a factor of 1,000-10,000 for dual and by 10,000-100,000 for triple redundancy. Redundancy essentially offers a trade-off between component quality and quantity. Our work in this paper permits this trade-off.

Another compelling study [7] uses an empirical assessment of how redundant computing improves time to solution. The simulation-driven study looked at a realistic scenario with a weak-scaling application that needs 168 hours to complete, a per-node MTTF of five years, a fixed five minutes to write out a checkpoint, and a fixed ten-minute time to restart. Checkpointing is performed at an optimal interval. The results show that at 200k nodes, an application spends eight times the time required to perform the work, reducing the throughput of such a machine to just over 10% compared to a fault-free environment. In contrast, using 400k nodes and dual redundancy, the elapsed wall clock time is 1/8 of that for the 200k-node non-redundant case. The throughput of the 400k-node system is four times better with redundant computing than the non-redundant 200k-node system. The prototype detailed in this paper is a step toward achieving this capability.

Other redundant MPI implementations, namely rMPI [33], MR-MPI [34], and VolpexMPI [35], do not protect against SDC. In contrast, RedMPI protects against SDC, lowers replication overheads, advances internal communications plus wildcard support, and provides SDC protected collectives that exploit native MPI collective performance.

Other parallel debugging tools that utilize redundancy exist such as MPI Echo [36][37].

VIII. CONCLUSION

Redundant computing is one approach to detect SDC. This study assesses the feasibility and effectiveness of SDC detection and correction at the MPI layer. We presented two

consistency protocols, explored the unique challenges in creating a deterministic MPI environment for replication purposes, investigated the effects of fault injection in to our framework, and analyzed the costs of performing SDC protection via redundancy.

This study develops a novel, efficient SDC detection and correction protocol (MsgPlusHash) with overheads ranging from 0% up to 30% for dual or triple redundancy depending on the number of messages sent by the application and the communication patterns. In particular, overheads do not change significantly for weak scaling applications as the number of processes is varied. These modest overhead ranges indicate the potential of RedMPI to protect against SDC for large-scale runs.

Our protocol detected and corrected injected faults for processes that continued to completion even when these faults resulted in many thousands of corrupted messages from a sender that experienced one or more SDC faults. In our controlled experiments, injected faults that were targeted to a single set of replicas were successfully isolated from spreading by fixing corrupted messages and allowing the applications to complete without incident. Further, when we injected faults into two or more replicas (of the same rank), RedMPI detected the corruption and was able to abort the application thus preventing invalid results from being reported.

In summary, RedMPI was successful in preventing invalid data from propagating or being transmitted without detection even under extreme scenarios. Our experiments showed profound effects from applications that experience even a single soft error without any form of protection. Without RedMPI, just one injected SDC was observed to quickly spread to other processes and messages, causing the majority of message data to become corrupt, which consistently lead to invalid results at a global scale.

Empirically, RedMPI not only performed exactly as expected in our controlled experiments, but it also pinpointed previously unknown hardware faults on our own experimental cluster nodes that had not been detected until RedMPI alerted us to SDCs occurring on unaltered MPI jobs.

While the cost of double/triple redundancy appears high in terms of power, analytic models show that for large core counts redundancy actually improves job throughput. As both the likelihood of node failure and silent data corruption increases as we scale up HPC systems, the importance of protecting data becomes obvious and available at a low cost when redundancy is already in place to ensure high throughput of mission-critical/high-consequence large scale applications.

ACKNOWLEDGEMENTS

Research sponsored by the Laboratory Directed Research and Development Program of Oak Ridge National Laboratory, managed by UT-Battelle, LLC for the U.S. Department of Energy under Contract No. De-AC05-00OR22725.

This research was partially supported by an Enterprise Partnership Scheme grant co-funded by IBM, the Irish Research Council for Science, Engineering & Technology (IRCSET), and the Industrial Development Agency (IDA) Ireland.

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 under contract DE-AC04-94AL85000.

This work was supported in part by NSF grants CNS-1058779, CNS-0958311.

REFERENCES

- [1] B. Schroeder, E. Pinheiro, and W.-D. Weber, "Dram errors in the wild: a large-scale field study," in *SIGMETRICS Conference on Measurement and Modeling of Computer Systems*, 2009, pp. 193–204.
- [2] E. Pinheiro, W.-D. Weber, and L. A. Barroso, "Failure trends in a large disk drive population," in *USENIX Conference on File and Storage Technologies*, 2007.
- [3] A. A. Hwang, I. A. Stefanovici, and B. Schroeder, "Cosmic rays don't strike twice: understanding the nature of dram errors and the implications for system design," in *Proceedings of the seventeenth international conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '12, 2012, pp. 111–122.
- [4] J. T. Daly, "ADTSC nuclear weapons highlights: Facilitating high-throughput ASC calculations," Los Alamos National Laboratory, Los Alamos, NM, USA, Tech. Rep. LALP-07-041, Jun. 2007.
- [5] J. T. Daly, L. A. Pritchett-Sheats, and S. E. Michalak, "Application MTTF vs. platform MTTF: A fresh perspective on system reliability and application throughput for computations at scale," in *Proceedings of the Workshop on Resiliency in High Performance Computing (Resilience) 2008*, May 2008, pp. 19–22.
- [6] I. Philp, "Software failures and the road to a petaflop machine," in *HPCRI: 1st Workshop on High Performance Computing Reliability Issues*, in *Proceedings of the 11th International Symposium on High Performance Computer Architecture (HPCA-11)*. IEEE Computer Society, 2005.
- [7] K. Ferreira, J. Stearley, J. H. L. III, R. Oldfield, K. Pedretti, R. Brightwell, R. Riesen, P. Bridges, and D. Arnold, "Evaluating the viability of process replication reliability for exascale systems," in *Supercomputing*, nov 2011.
- [8] A. Geist, "What is the monster in the closet?" Aug. 2011, invited Talk at Workshop on Architectures I: Exascale and Beyond: Gaps in Research, Gaps in our Thinking.
- [9] G. Bronevetsky and A. Moody, "Scalable i/o systems via node-local storage: Approaching 1 tb/sec file i/o," Lawrence Berkeley National Laboratory, TR 415791, 2009.
- [10] J. R. Sklaroff, "Redundancy management technique for space shuttle computers," *IBM Journal of Research and Development*, vol. 20, no. 1, pp. 20–28, 1976.
- [11] S. Mitra, N. Seifert, M. Zhang, Q. Shi, and K. S. Kim, "Robust system design with built-in soft-error resilience," *Computer*, vol. 38, no. 2, pp. 43–52, 2005.
- [12] M. Goma, C. Scarbrough, T. N. Vijayjumar, and I. Pomeranz, "Transient-fault recovery for chip multiprocessors," in *International Symposium on Computer Architecture*, May 2003, pp. 98–109.
- [13] S. K. Reinhardt and S. S. Mukherjee, "Transient fault detection via simultaneous multithreading," in *International Symposium on Computer Architecture*, 2000, pp. 25–36.
- [14] H. Quinn and P. Graham, "Terrestrial-based radiation upsets: A cautionary tale," in *Symposium on Field-Programmable Custom Computing Machines (FCCM) 2005*, Apr. 18–20, 2005, pp. 193–202.
- [15] J. Elliot, K. Kharbas, D. Fiala, F. Mueller, C. Engelmann, and K. Ferreira, "Combining partial redundancy and checkpointing for HPC," in *International Conference on Distributed Computing Systems*, 2012, p. (accepted).
- [16] J. Vetter, "Hpc landscape — application accelerators: Deus ex machina?" Sep. 2009, invited Talk at High Performance Embedded Computing Workshop.
- [17] J. Shalf, "Simulation challenge: Exascale planning overview," Aug. 2010, invited Talk at HEC FSIO R&D Workshop.
- [18] J. Dongarra, P. Beckman, T. Moore, P. Aerts, G. Aloisio, J. C. Andre, D. Barkai, J. Y. Berthou, T. Boku, B. Braunschweig, and et al., "The international exascale software project roadmap," *International Journal of High Performance Computing Applications*, vol. 25, no. 1, pp. 3–60, 2011.
- [19] D. Fiala, F. Mueller, C. Engelmann, K. Ferreira, R. Brightwell, and R. Riesen, "Detection and correction of silent data corruption for large-scale high-performance computing," Dept. of Computer Science, North Carolina State University, Tech. Rep. TR 2012-5, May 2012.
- [20] S. Böhm and C. Engelmann, "File i/o for mpi applications in redundant execution scenarios," in *EuroMicro International Conference on Parallel, Distributed, and network-based Processing*, Feb. 2012.
- [21] N. DeBardeleben, J. Laros, J. T. Daly, S. L. Scott, C. Engelmann, and B. Harrod, "High-end computing resilience: Analysis of issues facing the HEC community and path-forward for research and development," Whitepaper, Dec. 2009. [Online]. Available: <http://www.csm.ornl.gov/~engelman/publications/debardeleben09high-end.p%df>
- [22] P. H. Hargrove and J. C. Duell, "Berkeley Lab Checkpoint/Restart (BLCR) for Linux clusters," in *Journal of Physics: Proceedings of the Scientific Discovery through Advanced Computing Program (SciDAC) Conference 2006*, vol. 46. Denver, CO, USA: Institute of Physics Publishing, Bristol, UK, Jun. 25–29, 2006, pp. 494–499. [Online]. Available: http://www.iop.org/EJ/article/1742-6596/46/1/067/jpconf6_46_067.pdf
- [23] G. Bronevetsky and A. Moody, "Scalable I/O systems via node-local storage: Approaching 1 TB/sec file I/O," Lawrence Livermore National Laboratory, Livermore, CA, USA, Tech. Rep. TR-JLPC-09-01, Aug. 2009. [Online]. Available: <http://dx.doi.org/10.2172/964079>
- [24] S. Li, K. Chen, M.-Y. Hsieh, N. Muralimanohar, C. D. Kersey, J. B. Brockman, A. F. Rodrigues, and N. P. Jouppi, "System implications of memory reliability in exascale computing," in *Supercomputing*, 2011, pp. 46:1–46:12.
- [25] D. Fiala, K. Ferreira, F. Mueller, and C. Engelmann, "A tunable, software-based dram error detection and correction library for hpc," in *Workshop on Resiliency in High Performance Computing (Resilience) in Clusters, Clouds, and Grids*, Sep. 2011, pp. 110–121.
- [26] S. E. Michalak, K. W. Harris, N. W. Hengartner, B. E. Takala, and S. A. Wender, "Predicting the number of fatal soft errors in Los Alamos National Laboratory's ASC Q supercomputer," *IEEE Transactions on Device and Materials Reliability (TDMR)*, vol. 5, no. 3, pp. 329–335, 2005. [Online]. Available: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=1545893
- [27] G. Bronevetsky and B. R. de Supinski, "Soft error vulnerability of iterative linear algebra methods," in *Proceedings of the 21st ACM International Conference on Supercomputing (ICS) 2008*. Island of Kos, Greece: ACM Press, New York, NY, USA, Jun. 7–12, 2007. [Online]. Available: <http://greg.bronevetsky.com/papers/2008ICS.pdf>
- [28] D. P. Siemwiorek, "Architecture of fault-tolerant computers: An historical perspective," *Proceedings of the IEEE*, vol. 79, no. 12, pp. 1710–1734, 1991. [Online]. Available: <http://dx.doi.org/10.1109/5.119549>
- [29] A. Golander, S. Weiss, and R. Ronen, "DDMR: Dynamic and scalable dual modular redundancy with short validation intervals," *IEEE Computer Architecture Letters*, vol. 7, no. 2, pp. 65–68, 2008. [Online]. Available: <http://doi.ieeeecomputersociety.org/10.1109/L-CA.2008.12>
- [30] A. Shye, J. Blomstedt, T. Moseley, V. J. Reddi, and D. A. Connors, "PLR: A software approach to transient fault tolerance for multicore architectures," *IEEE Transactions on Dependable and Secure Computing (TDSC)*, vol. 6, no. 2, pp. 135–148, 2009. [Online]. Available: <http://doi.ieeeecomputersociety.org/10.1109/TDSC.2008.62>
- [31] S. S. Mukherjee, M. Kontz, and S. K. Reinhardt, "Detailed design and evaluation of redundant multithreading alternatives," in *Proceedings of the 29th Annual International Symposium on Computer Architecture (ISCA) 2002*. Anchorage, AK, USA: IEEE Computer Society, May 25–29, 2002, pp. 99–110. [Online]. Available: <http://doi.ieeeecomputersociety.org/10.1109/ISCA.2002.1003566>
- [32] C. Engelmann, H. H. Ong, and S. L. Scott, "The case for modular redundancy in large-scale high performance computing systems," in *Proceedings of the 8th IASTED International Conference on Parallel and Distributed Computing and Networks (PDCN) 2009*. Innsbruck, Austria: ACTA Press, Calgary, AB, Canada, Feb. 16–18, 2009, pp. 189–194. [Online]. Available: <http://www.csm.ornl.gov/~engelman/publications/engelmann09case.pdf>
- [33] R. Brightwell, K. B. Ferreira, and R. Riesen, "Transparent redundant computing with MPI," in *EuroMPI*, ser. Lecture Notes in Computer Science, R. Keller, E. Gabriel, M. M. Resch, and J. Dongarra, Eds., vol. 6305. Springer, 2010, pp. 208–218.
- [34] C. Engelmann and S. Böhm, "Redundant execution of hpc applications with mr-mpi," in *Proceedings of the 10th IASTED International Conference on Parallel and Distributed Computing and Networks (PDCN) 2011*. Innsbruck, Austria: ACTA Press, Calgary, AB, Canada, Feb. 15–17, 2011.
- [35] T. LeBlanc, R. Anand, E. Gabriel, and J. Subhlok, "Volpexmpi: An MPI library for execution of parallel applications on volatile nodes," in *Lecture Notes in Computer Science: Proceedings of the 16th European PVM/MPI Users' Group Meeting (EuroPVM/MPI)*

2009, vol. 5759. Espoo, Finland: Springer Verlag, Berlin, Germany, Sep. 7-10, 2009, pp. 124–133. [Online]. Available: http://dx.doi.org/10.1007/978-3-642-03770-2_19

- [36] B. Roundtree, G. Cobb, T. Gamblin, M. Schulz, B. Supinski, and H. Tufo, "Parallelizing heavyweight debugging tools with mpiecho," in *High-performance Infrastructure for Scalable Tools, WHIST 2011, Held as part of ICS '11, Tucson, Arizona*, 2011, pp. 803–808.
- [37] G. Cobb, B. Roundtree, H. Tufo, M. Schulz, T. Gamblin, and B. de Supinski, "Mpiecho: A framework for transparent mpi task replication," Dept. of Computer Science, University of Colorado at Boulder, Tech. Rep. CU-CS-1082-11, Jun. 2011.