

SANDIA REPORT

SAND2012-4060
Unlimited Release
Printed May, 2012

Evaluating Operating System Vulnerability to Memory Errors

Kurt B. Ferreira, Kevin Pedretti, Ron Brightwell, Patrick G. Bridges, David Fiala,
Frank Mueller

Prepared by
Sandia National Laboratories
Albuquerque, New Mexico 87185 and Livermore, California 94550

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.

Approved for public release; further dissemination unlimited.



Sandia National Laboratories

Issued by Sandia National Laboratories, operated for the United States Department of Energy by Sandia Corporation.

NOTICE: This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government, nor any agency thereof, nor any of their employees, nor any of their contractors, subcontractors, or their employees, make any warranty, express or implied, or assume any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represent that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government, any agency thereof, or any of their contractors or subcontractors. The views and opinions expressed herein do not necessarily state or reflect those of the United States Government, any agency thereof, or any of their contractors.

Printed in the United States of America. This report has been reproduced directly from the best available copy.

Available to DOE and DOE contractors from
U.S. Department of Energy
Office of Scientific and Technical Information
P.O. Box 62
Oak Ridge, TN 37831

Telephone: (865) 576-8401
Facsimile: (865) 576-5728
E-Mail: reports@adonis.osti.gov
Online ordering: <http://www.osti.gov/bridge>

Available to the public from
U.S. Department of Commerce
National Technical Information Service
5285 Port Royal Rd
Springfield, VA 22161

Telephone: (800) 553-6847
Facsimile: (703) 605-6900
E-Mail: orders@ntis.fedworld.gov
Online ordering: <http://www.ntis.gov/help/ordermethods.asp?loc=7-4-0#online>



Evaluating Operating System Vulnerability to Memory Errors

Kurt B. Ferreira (Org. 01423) kbferre@sandia.gov
Kevin Pedretti (Org. 01423) ktpedre@sandia.gov
Ron Brightwell (Org. 01423) rbbrigh@sandia.gov
Sandia National Laboratories
P. O. Box 5800
Albuquerque, NM 87185-1319

Patrick G. Bridges bridges@cs.unm.edu
Department of Computer Science
University of New Mexico

David Fiala fiala@cs.ncsu.edu
Frank Mueller mueller@cs.ncsu.edu
Department of Computer Science
North Carolina State University

Abstract

Reliability is of great concern to the scalability of extreme-scale systems. Of particular concern are soft errors in main memory, which are a leading cause of failures on current systems and are predicted to be the leading cause on future systems. While great effort has gone into designing algorithms and applications that can continue to make progress in the presence of these errors without restarting, the most critical software running on a node, the operating system (OS), is currently left relatively unprotected. OS resiliency is of particular importance because, though this software typically represents a small footprint of a compute

node's physical memory, recent studies show more memory errors in this region of memory than the remainder of the system. In this paper, we investigate the soft error vulnerability of two operating systems used in current and future high-performance computing systems: Kitten, the lightweight kernel developed at Sandia National Laboratories, and CLE, a high-performance Linux-based operating system developed by Cray. For each of these platforms, we outline major structures and subsystems that are vulnerable to soft errors and describe methods that could be used to reconstruct damaged state. Our results show the Kitten lightweight operating system may be an easier target to harden against memory errors due to its smaller memory footprint, largely deterministic state, and simpler system structure.

Introduction

Concern is growing in the high-performance computing (HPC) community on the reliability of future extreme scale systems. With systems continuing to grow dramatically in node count and individual nodes also increasing in component count and complexity, large-scale systems are becoming less reliable. In fact, experts are predicting that failure rates may go from the current state of a handful a day [35, 34] to multiple failures an hour [4]. Recent studies have shown soft errors in main memory to be the source of many of these failures [21, 27]. With the predicted increase of memory density on future exascale systems [38] and expected power optimizations such as decreases in supply voltages, the number of these failures is expected to dramatically increase.

Several methods have been developed to address these errors. Approaches include hardware-based techniques, such as single-bit error correction and double-bit detection (SEC-DED) and chipkill codes [12], as well as algorithm-based mechanisms that encode the correction mechanics directly into the application [20, 10, 7]. These hardware-based mechanisms may, however, be insufficient at the elevated failure rates predicted for exascale systems, and most importantly, they may not protect the most important software running on a node - the operating system.

An operating system (OS) resilient to soft errors in memory is key to the scalability of exascale systems for a number of reasons. First, current operating systems are unable to recover from the vast majority of failures. Second, though the typical operating system only occupies a small portion of a system's total physical memory footprint, recent studies show substantially more errors in this region than the remainder of a system's memory [21]. Lastly, future HPC system software will need to continue running in the presence of memory failures if current application-based, forward error recovery mechanisms are to be successful. These forward-error recovery methods are theorized to have lower overheads and less wasted computation than current rollback/recovery mechanisms.

In this work, we investigate the soft error vulnerability of two operating systems used in current and future high-performance computing systems: Kitten, the lightweight kernel developed at Sandia National Laboratories [33], and CLE, a high-performance HPC OS based on the Linux general purpose OS. Our analysis shows that the simpler lightweight kernel may be easier to harden against memory errors, because of its substantially smaller memory footprint, largely deterministic state, and generally simpler system structure.

Background

Current State of Practice

Coordinated checkpoint/restart is the dominant fault tolerance mechanism in high performance computing systems. In current systems, this approach works as follows:

1. Applications periodically quiesce all activity at a global synchronization point, for example a global barrier;
2. After synchronization, all nodes send some fraction of application and system state, generally comprising most of system memory, over the network to dedicated I/O nodes;
3. These I/O nodes store the received checkpoint information data to stable storage, currently hard disk-based storage;
4. In the event of application crash, the stored checkpoint can be used to restart the application at a prior known-good state.

The continued dominance of this technique rests on a number of key assumptions regarding failures that have thus far remained true:

1. Application state can be saved and restored much more quickly than a system's mean time to interrupt (MTTI);
2. The hardware and upkeep (e.g. power) costs of supporting frequent checkpointing are a modest portion (currently perhaps 10-20%) of the system's overall cost; and
3. System faults that do not crash (fail-stop) the system, such as so-called "soft errors", are very rare.

In an environment where failures are common, traditional checkpoint/restart has been shown to be inappropriate for large-scale systems [35, 2, 4, 18]. Additionally, checkpoint/restart is problematic when dealing with non-crash failures. In particular, checkpoint/restart *preserves* the impact of failures that corrupt application state. Addressing this problem requires application developers to either restart the application from scratch or analyze the contents of their checkpoints looking for one prior to when the fault that corrupted application state occurred.

Because of this limitation, there is significant effort underway within the community to develop *forward-error recovery* methods for application fault tolerance [18, 15, 6]. These methods deal with faults by correcting lost or incorrect state rather than restarting an application from a previously saved state. This approach avoids the wasted power and work of rollback/recovery methods like checkpointing and typically have significantly lower overheads.

DRAM Failures

Recent studies have shown DRAM errors in main memory to be the most common source of failures on today’s HPC platforms [21, 27]. The prevalence of these DRAM errors is related to the fact that typical large scale systems contain tens to hundreds of thousands of DRAM modules. A combination of the quantity and density of the information stored makes these modules particularly susceptible to faults. Moreover, with expected power optimizations, such as decreased supply voltages and increases in memory densities, the number of DRAM errors is expected to increase for future exascale systems [38].

To address these faults, current HPC systems typically include some form of error correction. The most common memory resilience scheme has the memory controller write additional checksum bits on each block of data. The memory controller then uses these checksum bits to detect and correct DRAM errors. Single-symbol Error Correction and Double-symbol Error Detection (SEC-DED) schemes allow systems to recover from the simplest memory failures and at least detect more complex (and less frequent) ones; or more complex chipkill-based codes [12] that allow a system to tolerate an entire DRAM chip failure at the cost of reduced performance and increased energy usage.

Uncorrectable DRAM errors, errors to two or more bits, are becoming increasingly common in systems with SEC-DED memory protection [36], with these errors occurring in up to 8% of DIMMs per year. For an exascale class system, this translates to multiple uncorrectable errors per hour. Such errors generally result in a machine check exception being delivered to the operating system, which then typically logs the error, and either kills the application to which the memory location belongs, or reboots the system if the error resides in a critical portion of the operating system’s address space [25].

As stated earlier, though the typical operating system occupies a very small portion of the system’s total physical memory, errors within the operating system’s address space are much more likely to occur than errors within the remainder of memory [21]. Therefore, techniques to address these errors at the system level are critical to the scalability of exascale systems.

Approach

The advantages described thus far in this paper provide a compelling reason to evaluate an HPC operating system’s vulnerability to memory errors. In this evaluation, we consider two operating systems of the type we expect to see on an exascale class system. The first is the Kitten lightweight kernel [33] developed by Sandia National Laboratories. The second is a variant of the Linux general-purpose operating system, called the Cray Linux environment.

Kitten is a special-purpose, limited-functionality OS designed for use on the compute nodes of massively parallel supercomputers. Its code base is derived from Linux, but is modified to minimize kernel-level functionality to only that needed for a set of mission-

critical HPC applications and moves as much as possible into user-space. Kitten is similar to previous lightweight kernels (LWK) such as SUNMOS, Puma, Cougar, and Catamount. Kitten, however, distinguishes itself from these prior LWKs by providing a Linux-compatible user environment, a more modern and extensible code base, and a virtual machine monitor capability via the Palacios virtual machine monitor [30] which allows full-featured guest operating systems to be loaded on-demand at very low overhead [26].

The Cray Linux Environment (CLE) is Cray’s scalable operating system for their XT line of supercomputers. CLE is based on the Linux general-purpose operating system with the addition of a number of optimizations to improve scalability. These optimizations include: enhancements to memory management, improved out-of-memory handling, and modifications for decreased OS jitter.

In this work, we will consider vulnerability to three types of common memory failures:

- Detected and corrected single-bit errors
- Detected but uncorrectable multi-bit errors
- Undetected “silent” data corruption

While fully protecting against each of these error types would be ideal, in many cases, the cost of doing so would far outweigh the benefit. Our goal is to identify the highest-impact opportunities for improving an OS’s resilience to memory errors.

Our evaluation will proceed as follows. First, for each OS, we will look at its complexity and how that complexity changes as a function of time. Our metric for complexity will be Source Lines Of Code (SLOC) count [11]. This metric gives us a rough measure of how difficult constructing and managing memory error mitigation methods will be. Next, we compare the memory footprints of the two operating systems, outlining how these footprints may change as an application progresses. Lastly, we breakdown the vulnerability of an OS on a per-subsystem basis, enumerating the subsystems’ critical state (state that must be free of errors). Additionally, for this critical state we describe possible failure mitigation strategies.

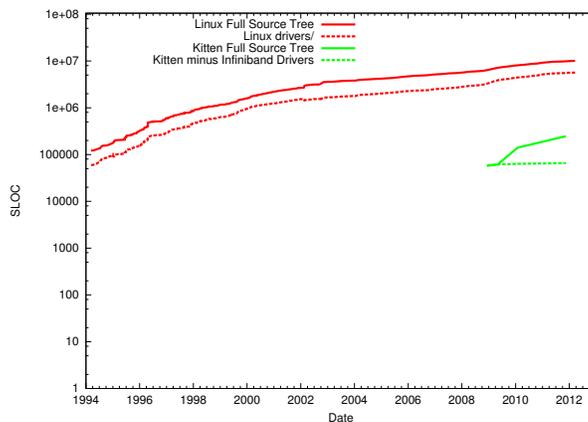
Results

Source Lines of Code

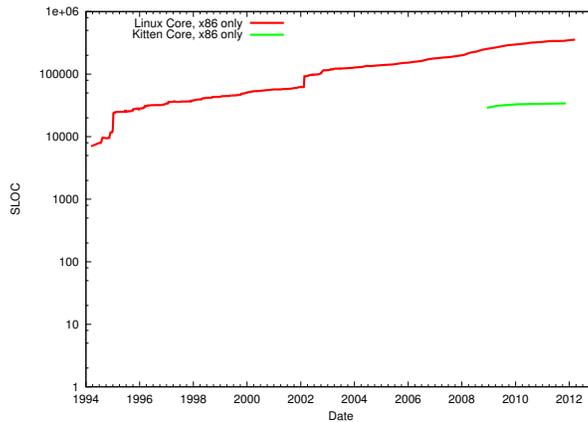
The Linux kernel has been enormously successful in attracting developers and users over its twenty year history. Due to this large development community and strong hardware support, Linux has also been successful in attracting HPC developers and is widely used within the community. Figure 1(a) plots the growth of the full Linux kernel codebase in

terms of source lines of code (SLOC), tracking its growth from approximately 120K SLOC in 1994 to its present size of over 10M SLOC. As the figure shows, the majority of the codebase consists of drivers. However, as shown in the right graph of Figure 1, non-driver core kernel code is also considerable and is growing rapidly. The current version of Linux, version 3.3, consists of approximately 350K SLOC in the core x86 architecture port (/kernel, /mm, and arch/x86 directories).

The Kitten codebase, in contrast, is currently a total of 246K SLOC, which drops to 66K SLOC once the Infiniband drivers and associated Linux driver support code are removed. Kitten’s core kernel code for the x86 architecture port is 30K SLOC, which is an order of magnitude smaller than the corresponding subset of Linux. This suggests that Kitten is considerably less complex than Linux, and will be easier to harden against memory errors.



(a) Full Tree



(b) Core Kernel Code, x86 Port Only

Figure 1. Comparison of Linux Kernel and Kitten Kernel source lines of code (SLOC).

Memory Footprint Comparison

Figure 2 compares the physical memory layouts used by Kitten and Linux. The primary difference between the two is that Kitten explicitly partitions memory into two regions, one for kernel memory and another for user-space applications, while Linux uses a unified page pool and dynamically assigns pages to different roles as needed. Kitten’s kernel memory footprint has a fixed upper limit (currently 64 MB) that does not change during runtime, while Linux’s footprint changes over time and can grow to the maximum size of physical memory.

Each user-space process on Kitten requires three pages of kernel memory to store task and address space structures, as well as a static amount of kernel memory to store the application’s page tables. When using 2 MB pages on the x86 architecture, approximately 8 KB of page table memory is needed for each gigabyte of application memory. Linux has similar per-process kernel memory requirements, with the addition of the pages in the page cache being used by the process. Kitten does not have a page cache. Additionally, Linux uses the 4 KB page size by default, resulting in more kernel memory being used for page tables (2 MB per GB of application memory). Libraries such as libhugetlbfs and recent transparent large page support in Linux are making it easier for applications to use large page sizes, with the caveat that memory fragmentation over time causes significant issues.

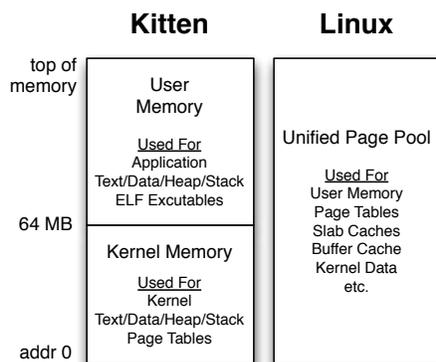


Figure 2. Physical memory layout of Kitten and Linux.

While a standard Linux kernel can grow to the full extent of the physical memory on the machine, typically the size is much smaller. In fact, CLE has a number of memory usage optimizations that limit memory footprint size. Specifically, CLE limits the size of the page cache using an I/O forwarding technique that avoids caching of file reads and writes. Figure 3 shows a comparison of the Kitten and CLE footprints. For Kitten, memory partitioning limits total kernel size to 64MB. For the CLE, we show the average kernel size measured using the smem [39] memory tracking tool while running the LAMMPS [32] molecular dynamics code from Sandia National Laboratories. From the figure, we see that

worst case Kitten OS size is more than an order of magnitude smaller than the average case from the CLE. Kitten’s smaller and deterministic footprint generally means simpler methods to protect and correct this state due to DRAM errors.

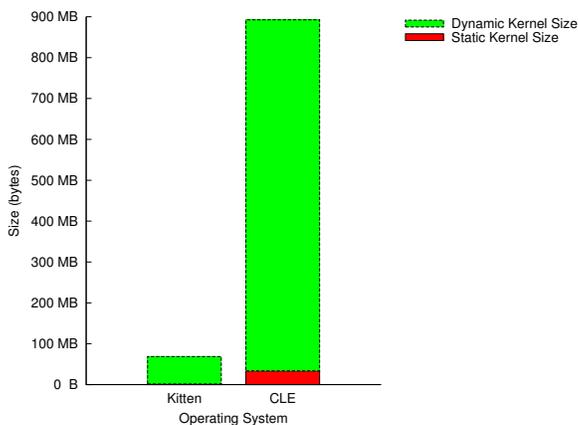


Figure 3. Comparison of the worst case Kitten static and dynamic kernel size to the average case measured on CLE. The average CLE memory footprint is an order of magnitude larger than the worst case for Kitten.

Major Kernel Subsystems

This section examines several kernel subsystems that exist in both operating systems, and discusses techniques that could be used to harden them against memory errors. The subsystems are discussed in the order of their kernel memory footprint in the Kitten kernel. This analysis captures the vast majority of Kitten’s kernel memory footprint, and is representative of the baseline kernel-level functionality needed to support highly-scalable HPC applications.

Page Table Memory

Both Kitten and Linux store page tables in kernel memory. The amount of page table memory used varies depending on the page size used: 4 KB pages require 2 MB of page table memory per gigabyte of application memory, 2 MB pages require 8 KB per gigabyte, and 1 GB pages require 8 bytes per gigabyte. In general, Kitten is always able to use the larger page sizes for application memory due to its segment-based and static memory allocation policy. Recent Linux kernels attempt to use large page sizes when possible, but memory fragmentation can limit the usefulness of this optimization.

Both OSs consider page table memory errors as fatal, either killing the affected application or the entire node. However, Kitten’s deterministic mapping of virtual to physical addresses would make it straightforward to recreate the corrupted page table memory contents from the base physical address and length information stored in the address space region object. This would not work on Linux due to its demand paging scheme, where unpredictable physical addresses are assigned to virtual addresses at runtime. Extra redundant state would need to be stored, and furthermore it may be difficult or impossible to tell which page table values have become corrupted if hardware notification is not provided.

Physical Memory Management

As described in Section , Linux uses a unified physical page pool. Linux maintains a memory map array, with one entry for each page of physical memory, to track the current state of each page frame in the system. Each entry in the table is 56 bytes, requiring 14 MB of overhead per GB of physical memory (1.4%).

Kitten does not maintain a memory map array. Instead, it maintains a free list of physical memory segments, where each segment consists of a physically contiguous set of pages with identical type (e.g., allocation status, memory type, associated NUMA node). Kitten’s segment list typically holds 10–100 entries for a high-end NUMA system, requiring less than 4 KB of kernel memory. Kitten’s physical memory tracking scheme would not work well for a general-purpose kernel, but is a good match for its target workloads where applications are allocated large contiguous regions of physical memory that are not demand paged.

As with page table memory, memory errors to the physical memory tracking data structures are considered fatal. In this case, however, there is no easy way to recreate the corrupted state. Instead, additional state of some form would have to be maintained, such as software-maintained ECC bits or redundant copies.

Dynamic Kernel Memory

Both OS kernels provide a mechanism for kernel subsystems and device drivers to allocate dynamic memory, similar to `malloc()` at user-level. Kitten implements this functionality using a buddy system memory allocator that covers the kernel memory portion of the physical address space (by default 64 MB). To avoid wasting memory due to over-allocation, Kitten uses a minimum block size of 32 bytes, which results in approximately 2 MB of buddy allocator state (one bit per 32 byte block). Kitten was profiled on an 8-core Intel system running an 8-thread OpenMP benchmark and found to use a maximum of 45 KB of dynamically allocated kernel memory.

Linux implements dynamic kernel memory allocation via a slab cache [5], which allocates physical memory from a buddy allocator that covers all of physical memory. The buddy

allocator uses a 4 KB minimum block size, resulting in approximately 256 KB of overhead per gigabyte of physical memory. Each slab cache requires a small state tracking structure of approximately 128 bytes plus 32 bytes per NUMA node. As an example, a Cray XE6 compute node running CLE 4.0.36 (Linux 2.6.32.45) maintains 130 slab caches of various sizes.

In addition to all of the memory allocator data structures being assumed to be reliable (buddy allocator state, slab cache info), each block of memory allocated has a small header at its start storing the size of the block and where it should be returned when freed (16 bytes for Kitten). This data would need to be protected from memory errors somehow, possibly by an ECC-like code or by storing redundant copies of the header. Alternatively, the caller could be made aware of the header so that it might try to protect it.

Address Spaces and Tasks

At its heart, Kitten's main purpose is to bootstrap user-space address spaces and tasks (threads and processes) and then get out of the way. Both address spaces and tasks are tracked by kernel-level state structures. Kitten's task structure is 8 KB in size and includes the task's kernel-level stack. Kitten's address space structure is around 800 bytes in size. Both structures are almost entirely self-contained, with only two pointers to additional data structures. Kitten's address space structure points to a list of virtual memory regions, of which there are usually four for a typical application address space: text, data, heap, and stack.

Linux has similar, but more complex task and address space structures. For example, the Linux task structure has over 160 fields, compared to 23 fields for Kitten. The obvious reason for this large difference is the additional functionality that Linux provides. However, much of this is not useful for HPC workloads, and simply increases the effort needed to understand and protect the codebase.

Kernel Entry and Exit

Kernel entry and exit occurs through well-defined interfaces. Both Linux and Kitten route all interrupts through a small assembly stub, which saves the necessary state and then calls the appropriate higher-level handler. Similarly, when applications make system calls, the kernel is entered through a common routine, which then redirects through a table to the appropriate handler.

This structure could potentially be leveraged to do coarse-grained kernel memory error detection and correction. At each kernel entry and exit, the entire kernel memory space could be checksummed to ensure that no kernel data was silently corrupted. This is straightforward on Kitten due to its contiguous kernel memory region. On Linux, kernel memory and application memory is interleaved both in physical memory and in the kernel's virtual address

space, making the checksum process more difficult but still possible.

Clearly, this approach would have high overhead when invoked. However, HPC applications typically do not make many system calls, and could benefit from the increased protection from memory errors. Additionally, it would eliminate the need to protect each individual kernel data structure, reducing memory overhead.

Page Retirement

An additional technique that applies to all of the kernel subsystems discussed thus far is page retirement [21]. In this scheme, the OS monitors the memory errors corrected by hardware and uses this information to predict which memory pages are likely to fail soon. Kernel data structures using these pages can then be migrated to more stable memory pages or discarded if appropriate. Recent versions of Linux can already use this technique to discard clean page cache pages that have experienced an uncorrectable memory error.

Kitten and Linux are both written in C, which makes migrating kernel data structures difficult since it is difficult to determine which other structures point to the data being moved. Furthermore, it is difficult to determine which kernel-level data structures are using a given page. Therefore, both OSs would require heavy modification in order to take advantage of this technique. In this regard, Kitten's smaller codebase could potentially be an advantage.

Related Work

Resiliency and fault-tolerance has been identified by the Department of Energy and Department of Defense as one of the key fundamental challenges of extreme-scale computing. The majority of the work in this active research area has focused solely on the application and ignored the operating and runtime systems, which is the focus of this work. Essentially all of these approaches attempt to improve the performance of checkpoint/restart as it is the most widely used mechanism for fault-tolerance today.

In addition to these application-based methods, a small handful of researchers have been focusing on designing fault-tolerant userspace libraries for HPC systems that applications can use to construct algorithm-based resilience. In each of these research areas is an underlying assumption that the operating and runtime systems are resilient to failures or if not, an expensive restart of the OS must be done. In the remainder of this section, we describe these approaches and briefly discuss their potential benefits and costs.

High-speed Storage for Checkpoint/Restart

Checkpointing to local disk and flash memory systems has periodically been proposed to speed up checkpoint/restart systems by placing large amounts of high-speed storage near the data that must be checkpointed. Actually deploying large amounts of local non-volatile storage in an exascale system is potentially very challenging. Local disk-based storage has traditionally been avoided because of the increased failures it causes, for example. Upcoming non-volatile phase change PCRAM, resistive RRAM devices, and modern NAND and NOR flash technologies provide high bandwidth and reliability, but are potentially very expensive. Unless their cost per bit rivals that of DRAM, using such technologies for checkpoint/restart purposes would result in checkpointing hardware that makes up a much larger portion of the system cost. Additionally, write durability issues may require periodically replacing all flash memory in the system, further impacting total costs.

Asynchronous Checkpointing and Message Logging

Another approach that has been suggested to improve the performance of checkpointing systems is uncoordinated or asynchronous checkpointing [1, 23, 24]. These methods typically checkpoint and restore from local storage without the synchronization used by coordinated checkpointing. To support a node restoring from a local asynchronous checkpoint, nodes in this approach keep a log of recent messages that they have sent. When a node restores from a previous checkpoint, it can then replay reception of messages using remote nodes' logs.

While this approach can increase checkpointing performance, logging increases the latency of messaging operations and potentially takes significant amounts of memory on a node. Finally, asynchronous checkpointing approaches can result in cascading rollbacks; recent work attempts to bound the amount of rollback that may be necessary [19], but also places non-trivial limits on application behavior. Lastly, thus far there has been little work examining the performance of a general message logging approach at the scales one might expect to see at exascale.

Other Checkpointing Systems

Memory-based checkpointing [31, 37] uses the memory of a remote machine to checkpoint node state. Unless node memory is primarily read-only (in which case RAID 5-like techniques can be used), this approach doubles the memory demands of an application. Since memory is regarded as a key budget and power constraint in exascale systems, the benefits of these techniques are unclear.

Multi-level checkpointing [29] is a library-based approach for controlling checkpointing to multiple storage targets, including memory-based checkpoints, local checkpoint storage, and remote checkpoints, into a single system. Because of this, it shares some of the advantages

and disadvantages of memory-based checkpointing and local storage techniques. Unlike these techniques, however, multi-level checkpointing has the flexibility to choose between multiple levels of storage based on system design parameters, making it a promising technique for exascale systems.

Finally, recent studies have looked at the benefits and costs of combining replication with traditional checkpoint/restart [18, 16, 14]. These studies seek to find the “break-even” points for replication, or the point where this replication approach uses fewer resources than traditional checkpoint/restart alone. In contrast to the other methods described thus far in this section, since replication typically duplicates not only the application processes but also a subset of the OS instances, errors with the operating and runtime system can be handled.

Fault Tolerant Userspace Libraries

In contrast to the checkpointing work described above, a number of researchers are investigating constructing libraries that are tolerant to certain kinds of faults. The idea being that the applications use these libraries to construct application-specific fault tolerance mechanisms, typically termed algorithm-based fault tolerance (ABFT) [20]. These ABFT techniques typically require a fault-tolerant message passing environment. There have been a number of these resilient message passing libraries based on MPI, including; FT-MPI [22, 17], AMPI [9], MPI/FT [3], and C^3 [8]. The differences between these libraries is beyond the scope of this work, but each of these libraries allows for an application to continue operating in the presence of faults, possibly in a degraded mode, and it is left up to the application to ensure the result is correct.

Current Operating System Memory Error Handling

OS-level handling of DRAM faults has generally been either very limited or used very heavyweight solutions. Linux and other operating systems, for example, provide low-level techniques for handling, logging, and notifying the application of such errors [25]. These techniques generally terminate the application or OS kernel, and potentially invoke higher-level recovery systems based on, for example, checkpointing or redundancy. Some systems have attempted to provide additional protection against memory faults both on CPUs [13] and GPUs [28], though with substantial cost.

Conclusions and Future Work

In this paper we presented a preliminary evaluation of operating system vulnerability to DRAM failures, a common error in current and future extreme-scale systems. Hardening system software to this class of errors will be critical for the success of emerging fault-tolerance

methods. This work focused on two HPC operating systems; Kitten, the lightweight operating system developed at Sandia National Laboratories and the Cray Linux Environment, a HPC variant of the Linux operating system. Each of these OSs represents an OS construction methodology currently used in HPC. For each OS, we present the complexity of each OS in terms of the metric SLOCCount, examine the memory footprint, and evaluate vulnerability on a per subsystem basis. Where critical state is found, state that must be protected from DRAM errors, we outline mitigation methods that can be used. Overall, these results suggest hardening the Kitten lightweight kernel to be more tractable due to its smaller and deterministic state in comparison to Linux.

While this preliminary analysis shows there is promise in this idea, more work is clearly needed. For example, a detailed analysis of the mitigation techniques, outlining both the space and performance overheads is needed to decide which methods are ideal. Additionally, hardening system software to failures beyond those that occur in system RAM will be key to scalability of extreme-scale systems. Lastly, evaluating the hardened OSs and system services to errors will be key to outlining this work's overall merit.

References

- [1] Jinho Ahn. 2-step algorithm for enhancing effectiveness of sender-based message logging. In *SpringSim '07: Proceedings of the 2007 spring simulation multiconference*, pages 429–434, 2007.
- [2] Saman Amarasinghe and et al. Exascale software study: Software challenges in extreme scale systems. <http://users.ece.gatech.edu/mrichard/ExascaleComputingStudyReports/ECSS%20report%20101909.pdf>, September 2009.
- [3] Rajanikanth Batchu, Yiogonder S. Dandass, Anthony Skjellum, and Murali Beddhu. MPI/FT: A model-based approach to low-overhead fault tolerant message-passing middleware. *Cluster Computing*, 7(4):303–315, Jan. 2004.
- [4] Keren Bergman, Shekhar Borkar, Dan Campbell, William Carlson, William Dally, Monty Denneau, Paul Franzon, William Harrod, Kerry Hill, Jon Hiller, Sherman Karp, Stephen Keckler, Dean Klein, Peter Kogge, Robert Lucas, Mark Richards, Al Scarpelli, Steven Scott, Allan Snavely, Thomas Sterling, R. Stanley Williams, and Katherine Yelick. Exascale computing study: Technology challenges in achieving exascale systems. [http://www.science.energy.gov/ascr/Research/CS/DARPAexascale-hardware\(2008\).pdf](http://www.science.energy.gov/ascr/Research/CS/DARPAexascale-hardware(2008).pdf), September 2008.
- [5] Jeff Bonwick and Jonathan Adams. Magazines and vmem: Extending the slab allocator to many CPUs and arbitrary resources. In *Proceedings of the General Track: 2002 USENIX Annual Technical Conference*, pages 15–33, Berkeley, CA, USA, 2001. USENIX Association.
- [6] Patrick Bridges, Mark Hoemmen, Kurt B Ferreira, Michael Heroux, Philip Soltero, and Ron Brightwell. Cooperative application/os DRAM fault recovery. *Workshop on Resiliency in High Performance Computing (Resilience) in Clusters, Clouds, and Grids in conjunction with the Euro-Par Conference, Lecture Notes in Computer Science*, pages –, 2011.
- [7] Greg Bronevetsky and Bronis de Supinski. Soft error vulnerability of iterative linear algebra methods. In *Proceedings of the 22nd Annual International Conference on Supercomputing*, ICS '08, pages 155–164, New York, NY, USA, 2008. ACM.
- [8] Greg Bronevetsky, Daniel Marques, Keshav Pingali, and Paul Stodghill. Collective operations in application-level fault-tolerant MPI. In *Proceedings of the 17th annual international conference on Supercomputing*, ICS '03, pages 234–243, New York, NY, USA, 2003. ACM.
- [9] Sayantan Chakravorty, Celso Mendes, and Laxmikant Kal. Proactive fault tolerance in mpi applications via task migration. *Strategy*, 4297:485496, 2006.

- [10] Zizhong Chen and J. Dongarra. Algorithm-based checkpoint-free fault tolerance for parallel matrix computations on volatile resources. In *Parallel and Distributed Processing Symposium, 2006. IPDPS 2006. 20th International*, April 2006.
- [11] David A. Wheeler. Sloccount. <http://www.dwheeler.com/sloccount>, March 1 2012.
- [12] Timothy J. Dell. A white paper on the benefits of chipkill-correct ECC for PC server main memory. IBM Microelectronics Division, Nov. 1997.
- [13] Dave Dopson. SoftECC: A system for software memory integrity checking. Master's thesis, Massachusetts Institute of Technology, September 2005.
- [14] James Elliot, Kishor Kharbas, David Fiala, Frank Mueller, Kurt Ferreira, and Christian Engelmann. Combining partial redundancy and checkpointing for HPC. In *International Conference on Distributed Computing Systems*, pages 1–11, Los Alamitos, CA, USA, June 2012. IEEE Computer Society Press. [to appear].
- [15] Christian Engelmann and George A. (Al) Geist. Super-scalable algorithms for computing on 100,000 processors. In *Lecture Notes in Computer Science: Proceedings of the 5th International Conference on Computational Science (ICCS) 2005, Part I*, volume 3514, pages 313–320, Atlanta, GA, USA, May 22-25, 2005. Springer Verlag, Berlin, Germany.
- [16] Christian Engelmann, Hong H. Ong, and Stephen 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*, pages 189–194, Innsbruck, Austria, February 16-18, 2009. ACTA Press, Calgary, AB, Canada.
- [17] Graham E. Fagg, Thara Angskun, George Bosilca, Jelena Pjesivac-Grbovic, and Jack Dongarra. Scalable fault tolerant mpi: Extending the recovery algorithm. In Beniamino Di Martino, Dieter Kranzlmüller, and Jack Dongarra, editors, *PVM/MPI*, volume 3666 of *Lecture Notes in Computer Science*, pages 67–75. Springer, 2005.
- [18] Kurt Ferreira, Rolf Riesen, Jon Stearley, James H. Laros III, Ron Oldfield, Kevin Pedretti, Patrick Bridges, Dorian Arnold, and Ron Brightwell. Evaluating the viability of process replication reliability for exascale systems. In *Proceedings of the ACM/IEEE International Conference on High Performance Computing, Networking, Storage, and Analysis, (SC'11)*, Nov 2011.
- [19] Amina Guermouche, Thomas Ropars, Elisabeth Brunet, Marc Snir, and Franck Cappello. Uncoordinated checkpointing without domino effect for send-deterministic message passing applications. In *Proceedings of the 2011 IEEE International Parallel and Distributed Processing Symposium*, May 2011.
- [20] Kuang-Hua Huang and Jacob A. Abraham. Algorithm-based fault tolerance for matrix operations. *IEEE Transactions on Computers*, C-33(6), June 1984.

- [21] Andy A. Hwang, Ioan A. Stefanovici, and Bianca 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, ASPLOS '12*, pages 111–122, New York, NY, USA, 2012. ACM.
- [22] Inovative Computing Laboratory. FT-MPI. <http://icl.cs.utk.edu/ftmpi>, March 1 2012.
- [23] Q. Jiang and D. Manivannan. An optimistic checkpointing and selective approach for consistent global checkpoint collection in distributed systems. In *Proceedings of the 2007 IEEE International Parallel and Distributed Processing Symposium*, March 2007.
- [24] David B. Johnson and Willy Zwaenepoel. Recovery in distributed systems using asynchronous and checkpointing. In *Proceedings of the seventh annual ACM Symposium on Principles of distributed computing*, pages 171–181, 1988.
- [25] Andi Kleen. mcelog: memory error handling in user space. In *Proceedings of Linux Kongress 2010*, Nuremburg, Germany, September 2010.
- [26] John R. Lange, Kevin T. Pedretti, Trammell Hudson, Peter A. Dinda, Zheng Cui, Lei Xia, Patrick G. Bridges, Andy Gocke, Steven Jaconette, Michael Levenhagen, and Ron Brightwell. Palacios and kitten: New high performance operating systems for scalable virtualized and native supercomputing. In *IPDPS'10*, pages 1–12, 2010.
- [27] Sheng Li, Ke Chen, Ming-Yu Hsieh, Naveen Muralimanohar, Chad D. Kersey, Jay B. Brockman, Arun F. Rodrigues, and Norman P. Jouppi. System implications of memory reliability in exascale computing. In *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis, SC '11*, pages 46:1–46:12, New York, NY, USA, 2011. ACM.
- [28] N. Maruyama, A. Nukada, and S. Matsuoka. A high-performance fault-tolerant software framework for memory on commodity GPUs. In *Parallel Distributed Processing (IPDPS), 2010 IEEE International Symposium on*, pages 1 –12, april 2010.
- [29] Adam Moody, Greg Bronevetsky, Kathryn Mohror, and Bronis R. de Supinski. Design, modeling, and evaluation of a scalable multi-level checkpointing system. In *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis, SC '10*, pages 1–11, Washington, DC, USA, 2010. IEEE Computer Society.
- [30] Northwestern University. Palacios: An os independent embeddable vmm. <http://v3vee.org/palacios>, March 10 2012.
- [31] J. S. Plank, Y. B. Kim, and J. J. Dongarra. Algorithm-based diskless checkpointing for fault tolerant matrix operations. In *Twenty-Fifth International Symposium on Fault-Tolerant Computing. Digest of Papers*, pages 351–360, Pasadena, CA, USA, June 1995. Los Alamitos, CA, USA : IEEE Comput. Soc. Press, 1995.

- [32] Sandia National Laboratories. The LAMMPS molecular dynamics simulator. <http://lammps.sandia.gov>, April 2010.
- [33] Sandia National Laboratory. Kitten lightweight kernel. <https://software.sandia.gov/trac/kitten>, March 10 2012.
- [34] Bianca Schroeder and Garth A. Gibson. A large-scale study of failures in high-performance computing systems. In *Proceedings of the International Conference on Dependable Systems and Networks (DSN2006)*, June 2006.
- [35] Bianca Schroeder and Garth A Gibson. Understanding failures in petascale computers. *Journal of Physics: Conference Series*, 78(1):012022, 2007.
- [36] Bianca Schroeder, Eduardo Pinheiro, and Wolf-Dietrich Weber. DRAM errors in the wild: a large-scale field study. *Communications of the ACM*, 54:100–107, February 2011.
- [37] L. M. Silva and J. G. Silva. An experimental study about diskless checkpointing. In *24th EUROMICRO Conference*, pages 395 – 402, Vasteras, Sweden, August 1998. IEEE Computer Society Press.
- [38] Horst Simon. Exascale challenges for the computational science community. Technical report, Lawrence Berkeley National Laboratory and UC Berkeley, Oct. 2010.
- [39] SMEM. Memory reporting tool. <http://www.selenic.com/smем/>, March 1 2012.

DISTRIBUTION:

1	MS 1319	Kurt Ferreira , 1423
1	MS 1319	Kevin Pedretti , 1423
1	MS 1319	Ron Brightwell , 1423
1	MS 0899	Technical Library, 9536 (electronic copy)
1	MS 0359	D. Chavez, LDRD Office, 1911

