

Heterogeneous-Reliability Memory: Exploiting Application-Level Memory Error Tolerance

Yixin Luo Sriram Govindan[†] Bikash Sharma[†] Mark Santaniello[†] Justin Meza
Aman Kansal[†] Jie Liu[†] Badriddine Khessib[†] Kushagra Vaid[†] Onur Mutlu
Carnegie Mellon University, yixinluo@cs.cmu.edu, {meza, onur}@cmu.edu

[†]Microsoft Corporation, {srgovin, bsharma, marksan, kansal, jie.liu, bk Hessib, kvoid}@microsoft.com

1. Summary

Recent studies estimate that server cost contributes to as much as 57% of the total cost of ownership (TCO) of a data-center [1]. One key contributor to this high server cost is the procurement of memory devices such as DRAMs, especially for data-intensive datacenter cloud applications that need low latency (such as web search, in-memory caching, and graph traversal). Such memory devices, however, may be prone to hardware errors that occur due to unintended bit flips during device operation [40, 33, 41, 20]. To protect against such errors, traditional systems uniformly employ devices with high-quality chips and error correction techniques, both of which increase device cost. At the same time, we make the observations that 1) data-intensive applications exhibit a diverse spectrum of tolerance to memory errors, and 2) traditional one-size-fits-all memory reliability techniques are *inefficient* in terms of cost.

Our DSN-44 paper [30] is the first to 1) understand how tolerant different data-intensive applications are to memory errors and 2) design a new memory system organization that matches hardware reliability to application tolerance in order to reduce system cost. The **main idea** of our approach is to classify applications based on their memory error tolerance, and map applications to *heterogeneous-reliability* memory system designs managed cooperatively between hardware and software to reduce system cost. Our DSN-44 paper provides the following **contributions**:

1. A new **methodology** to quantify the tolerance of applications to memory errors. Our approach measures the effect of memory errors on application correctness and quantifies an application’s ability to mask or recover from memory errors.
2. A comprehensive **characterization** of the memory error tolerance of three data-intensive workloads: an interactive web search application [30, 39], an in-memory key–value store [30, 3], and a graph mining framework [30, 29]. We find that there exists an order of magnitude difference in memory error tolerance across these three applications.
3. An **exploration** of the design space of new memory system organizations, called *heterogeneous-reliability memory*, which combines a heterogeneous mix of reliability techniques that leverage application error tolerance to reduce system cost. We show that our techniques can reduce server hardware cost by 4.7%, while achieving 99.90% single server availability.

1.1. Methodology

We had three design goals when implementing our methodology for quantifying application memory error tolerance. First, due to the sporadic and inconsistent nature of memory

errors in the field [40, 33, 41, 20, 27, 18, 36], we wanted to design a framework to emulate the occurrence of a memory error in an application’s data in a *controlled* manner. Second, we wanted an *efficient* way to measure how an application accesses its data. Third, we wanted our framework to be easily *adaptable* to other workloads or system configurations.

Figure 1 shows a flow diagram illustrating the five steps involved in our error emulation framework. We assume that the application under examination has already been run outside of the framework and its expected output has been recorded. The framework proceeds as follows. (1) We start the application under the error injection framework.¹ (2) We use software debuggers² to inject the desired number and types of memory errors. (3) We initiate the connection of a client and start executing the desired workload. (4) Throughout the course of the application’s execution, we check to see if the machine has crashed; if it has, we log this outcome and proceed to step (1) to begin testing once again. (5) If the application finishes its workload, we check to see if its output matches the expected results; we log this outcome and proceed to step (1) to test again.

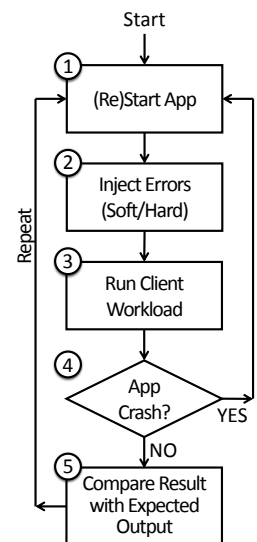


Figure 1. Memory error emulation framework.

1.2. Characterization

In our DSN paper, we characterize three new data-intensive applications to memory errors, to quantify their tolerance to memory errors: **WebSearch** [39]—an interactive web search application, **Memcached** [3]—an in-memory key-value store, and **GraphLab** [29]—a graph mining framework. We ran these applications in real production systems and sampled hundreds to tens of thousands of unique memory addresses for each application.

Key Findings. In the following, we summarize two of the most important findings of our DSN 2014 paper.

Finding 1: Error Tolerance Varies Across Applications. Figure 2(a) plots the probability of each of the three applications crashing due to the occurrence of single-bit soft or hard

¹We describe our memory error emulation framework in our DSN 2014 paper [30].

²WinDBG [4] in Windows and GDB [2] in Linux.

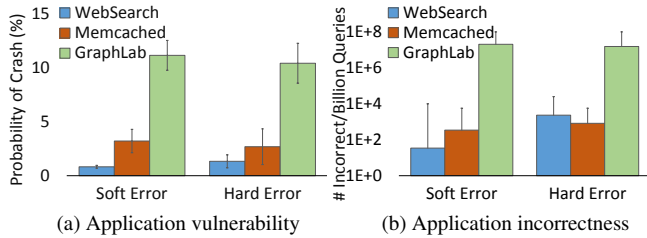


Figure 2. Inter-application variations in vulnerability to single-bit soft and hard memory errors for the three applications in terms of (a) probability of crash and (b) frequency of incorrect results.

errors in their memory (we call this application-level memory error vulnerability). In terms of how these errors affect application correctness, Figure 2(b) plots the rate of incorrect results per billion queries under the same conditions. We draw two key observations from these results.

First, there exists a significant variance in vulnerability among the three applications both in terms of crash probability and in terms of incorrect result rate, which varies by up to six orders of magnitude. Second, these characteristics may differ depending on whether errors are soft or hard (for example, the number of incorrect results for WebSearch differs by over two orders of magnitude between soft and hard errors). We therefore conclude that *memory reliability techniques that treat all applications similarly are inefficient because there exists significant variation in error tolerance among applications.*

Finding 2: Error Tolerance Varies Within an Application.

Figure 3(a) plots the probability of each of the three applications crashing due to the occurrence of single-bit soft or hard errors in *different regions* of their memory. Figure 3(b) plots the rate of incorrect results per billion queries under the same conditions. Two observations are in order.

First, for some memory regions, the probability of an error leading to a crash is much lower than for others (for example, in WebSearch the probability of a hard error leading to a crash in the *heap* or *private* memory regions is much lower than in the *stack* memory region). Second, even in the presence of memory errors, some regions of some applications are still able to tolerate memory errors (perhaps at reduced correctness). This may be acceptable for applications such as WebSearch that aggregate results from several servers before presenting them to the user, in which case the likelihood of the user being exposed to an error is much lower than the reported probabilities. We therefore conclude that *memory reliability techniques that treat all memory regions within an application similarly are inefficient because there exists significant variance in the error tolerance among different memory regions.*

Other Findings. We use WebSearch as an exemplar data-intensive application and perform an in-depth analysis of its tolerance behavior. Other findings in the paper are listed below.

- Quick-to-crash vs. periodically incorrect behavior. More-severe failures due to memory errors tend to crash the application or system quickly, while less-severe failures tend to generate incorrect results periodically.

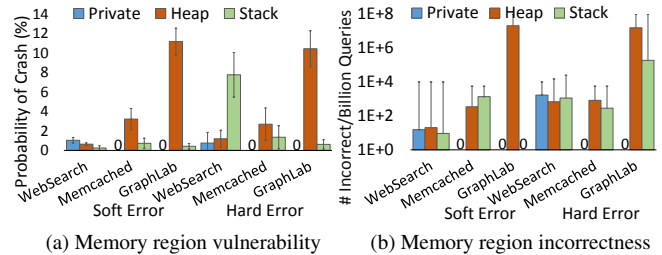


Figure 3. Memory region variations in vulnerability to single-bit soft and hard memory errors for the applications in terms of (a) probability of crash and (b) frequency of incorrect results.

- Some memory regions are safer than others. This indicates either application access pattern or application access pattern or application logic can be the dominant factor in different memory regions to mask a majority of memory errors.
- More-severe errors mainly decrease correctness, as opposed to increase an application’s probability of crashing.
- Data recoverability varies across memory regions. For data-intensive applications like WebSearch, software-only memory error tolerance techniques are a promising direction for enabling reliable system designs.

1.3. Solution

In our DSN paper, we propose *heterogeneous-reliability memory* as our solution to optimize datacenter cost. We examine three dimensions and their benefits and trade-offs in the design space for systems with heterogeneous reliability: (1) hardware techniques to detect and correct errors, (2) software responses to errors, and (3) the granularity at which different techniques are used. In addition to dimensions in system designs, we also discuss (1) the metrics we use to evaluate the benefits and costs of the designs, and (2) the memory error model we use to examine the effectiveness of the designs.

Using WebSearch as an example application, we evaluated and compared five design points (three baseline systems: Typical Server, Consumer PC, and Less-Tested, and two heterogeneous-reliability memory systems: Detect&Recover and Detect&Recover/L). Such comparison illustrates the inefficiencies of traditional homogeneous approaches to memory system reliability, as well as the benefits of heterogeneous-reliability memory system designs. Detect&Recover uses parity in hardware to detect errors for the *private* memory region and responds by correcting them with a clean copy of data from disk in software (Par+R, parity and recovery), and uses neither error detection nor correction for the rest of its data. Detect&Recover/L applies the same technique as Detect&Recover does but uses less-thoroughly-tested memory throughout the entire system.

Major evaluation results in Figure 4 show that the two highlighted design points (in orange color), which leverage *heterogeneous-reliability* memory system design, can both achieve the target single server availability of 99.90% while reducing server hardware cost to 2.9% and 4.7% respectively. We therefore conclude that *heterogeneous-reliability memory system designs can enable systems to achieve both high cost*

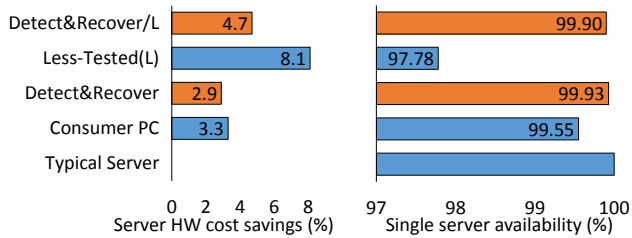


Figure 4. Comparison of server hardware cost savings and single server availability for the five design points.

savings and high single server availability/reliability at the same time.

Our DSN 2014 paper has much more detailed analysis, including 1) memory cost savings, expected crash and incorrect query frequency for each configuration (Sec. VI-B), 2) the maximum number of tolerable errors per month for each application to achieve a reliability target (Sec. VI-B), and 3) a discussion of hardware/software support for and feasibility of the proposed heterogeneous-reliability memory (Sec. VI-C).

2. Significance

2.1. Novelty

Our DSN paper identifies the inefficiency in traditional *one-size-fits-all* solutions to tolerate memory errors. We find an opportunity to greatly reduce server hardware cost by provisioning the right amount of memory reliability for different applications using a hardware-software cooperative approach. By characterizing three important data-intensive applications, we build a comprehensive understanding of how tolerant different data-intensive applications are to memory errors. Based upon these insights, we propose the idea of heterogeneous-reliability memory system and explore dimensions in the design space of such systems.

Methodology. To perform a controlled memory error injection to the three new data-intensive applications in a rapid manner, we build a unique memory error emulation framework using software debuggers [4, 2]. This framework allows us to meet our goal of efficiency and adaptivity coming from the scale and diversity of our target applications and enables us to execute our workloads in their entirety. This framework enables us to understand different aspects of application memory error tolerance at the whole-application level.

Characterization. We characterize different aspects of an application’s tolerance to memory errors, from the resulting outcomes to the underlying causes. We use our unique characterization framework to inject a total of 23,718 unique memory errors and monitor the correctness of a total of 11 billion query results for the three new data-intensive applications. From the characterization results, we conclude 6 new findings which motivate and enable our proposed solution.

Solution. To reduce datacenter cost while reaching a reliability target, we propose a heterogeneous-reliability memory system design, managed cooperatively between hardware and software. We examine three dimensions in the design space enabled by this system design. Our evaluations of an example heterogeneous-reliability memory system show that our proposed solution can enable future low-cost high-reliability memory system designs with minimal changes to hardware and software.

2.2. Long-Term Impact

We believe our DSN paper will have long-term impact due to the following three reasons. First, it emphasizes and aims to solve an unfavorable trend in DRAM scaling — an increase in hardware cost to ensure memory reliability. Second, it tackles a problem that will become increasingly important in the future — reducing memory system cost in datacenters. Third, it proposes a technique that uses software-hardware co-design to improve memory system reliability, thereby hopefully inspiring future works to exploit software characteristics to improve system reliability and reduce system cost.

Increasing memory error rate. As DRAM scales to smaller process technology nodes, DRAM reliability degradation becomes more significant [17, 34]. For example, recent works 1) showed the existence of disturbance errors in commodity DRAM chips operating in the field [21], 2) experimentally demonstrated the increasing importance of retention related failures in modern DRAM devices [26, 28, 19, 36], and 3) advocated, in a paper co-written by Samsung and Intel design teams, the use of in-DRAM error correcting codes to overcome the reliability challenges [17, 36]. As a result of decreasing DRAM reliability, keeping the memory error rate as the levels we have today can either increase DRAM cost due to decreased yield, expensive quality assurance tests, and extra capacity for storing stronger error-correcting codes or reduce performance due to frequent error correction and logging. All of these make DRAM scaling more difficult and less appealing. Our paper proposes a solution that enables the use of DRAMs with higher error rates while achieving reasonable application reliability, which can enable more efficient scaling of DRAM to smaller technology nodes in the future.

Other memory technologies such as NAND Flash Memory [12, 11, 10, 9, 6, 7, 8, 32], phase change memory (PCM) [23, 24, 23, 37, 31] and STT-MRAM [22, 31] also show a similar decreasing trend in their reliability with process technology scaling and the advent of multi-level cell (MLC) technology [35]. NAND Flash Memory, for example, also suffers from retention errors [6, 7, 12] and cell-to-cell program interference errors [6, 9, 11]. Additionally, NAND Flash Memory suffers from program/erase cycling errors [6, 10], and read disturb errors [6, 13]. PCM suffers from endurance issues [23, 37, 38, 25] and resistance drift [16]. The solution proposed in our paper can also be applied to these memory technologies with slight modification to enable reliable high density non-volatile devices in the future.

Increasing datacenter cost. Recent studies have shown that capital costs can account for the majority (e.g., around 57% in [5]) of datacenter TCO (total cost of ownership). As part of the cost of a server, the cost of the memory is comparable to that of the processors, and is likely to exceed processor cost and become the dominant cost for servers running data-intensive applications such as web search and social media services. As future datacenters grow in scale, datacenter TCO will become an increasingly important factor in system design. Our paper demonstrates a way of optimizing datacenter TCO by reducing the cost of the memory system. The cost savings can be significant due to the increasing scale of such datacenters [15], making our proposed technique hopefully

more important in the future.

Software-hardware co-design. Our solution, heterogeneous-reliability memory, utilizes software-hardware cooperative design to reduce cost. It demonstrates the benefits of exploiting the application characteristics to improve overall system design. For example, it shows that a significant number of errors can be corrected in software by reloading a clean copy of the data from storage. This motivates us to rethink the placement of different functionalities (such as error detection and error correction) to improve the cost-reliability trade-off.

2.3. New Research Directions

Our DSN paper will hopefully continue to inspire future works that can provide efficient and extensive characterization/estimation of application-level memory error tolerance [14], which can make our proposed technique applicable to a broader set of applications. Our paper has started a community discussion [15] on the feasibility of solving the problem of memory reliability by exploiting application memory error tolerance in the future, inspiring reporters to ask the question: "How good does memory need to be?"

2.4. Conclusion

Our DSN paper presents an efficient methodology to characterize memory error tolerance, a set of new characterization results along with their findings, and a new solution to reduce datacenter cost under an availability constraint. As DRAM becomes increasingly less reliable at smaller technology nodes and memory cost becomes more important in datacenters in the future, we hope that our findings and ideas will inspire more research to improve the cost-reliability trade-off in memory systems.

References

- [1] Data Center Cost Model. <http://perspectives.mvdirona.com/2010/09/18/OverallDataCenterCosts.aspx>.
- [2] GDB: The GNU Project Debugger. <http://www.sourceware.org/gdb/>.
- [3] Memcached. <http://memcached.org/>.
- [4] Windows Debugging. <http://tinyurl.com/l6zsqzv>.
- [5] L. A. Barroso et al. *The Datacenter as a Computer: An Introduction to the Design of Warehouse-Scale Machines*. Morgan & Claypool Publishers, 2009.
- [6] Y. Cai et al. Error Patterns in MLC NAND Flash Memory: Measurement, Characterization, and Analysis. In *DATE*, 2012.
- [7] Y. Cai et al. Flash Correct-and-Refresh: Retention-Aware Error Management for Increased Flash Memory Lifetime. In *ICCD*, 2012.
- [8] Y. Cai et al. Error Analysis and Retention-Aware Error Management for NAND Flash Memory. *ITJ*, 2013.
- [9] Y. Cai et al. Program Interference in MLC NAND Flash Memory: Characterization, Modeling, and Mitigation. In *ICCD*, 2013.
- [10] Y. Cai et al. Threshold Voltage Distribution in MLC NAND Flash Memory: Characterization, Analysis, and Modeling. In *DATE*, 2013.
- [11] Y. Cai et al. Neighbor-Cell Assisted Error Correction for MLC NAND Flash Memories. In *SIGMETRICS*, 2014.
- [12] Y. Cai et al. Data Retention in MLC NAND Flash Memory: Characterization, Optimization and Recovery. In *HPCA*, 2015.
- [13] Y. Cai et al. Read Disturb Errors in MLC NAND Flash Memory: Characterization, Mitigation and Recovery. In *DSN*, 2015.
- [14] N. Foutris et al. Versatile Architecture-level Fault Injection Framework for Reliability Evaluation: A first Report. In *IOLTS*, 2014.
- [15] R. Harris. How good does memory need to be?, 2014. <http://www.zdnet.com/how-good-does-memory-need-to-be-7000031853/>.
- [16] D. Ielmini et al. Physical Interpretation, Modeling and Impact on Phase Change Memory (PCM) Reliability of Resistance Drift Due to Chalcogenide Structural Relaxation. In *IEDM*, 2007.
- [17] U. Kang et al. Co-Architecting Controllers and DRAM to Enhance DRAM Process Scaling. In *The Memory Forum*, 2014.
- [18] S. Khan et al. The Efficacy of Error Mitigation Techniques for DRAM Retention Failures: A Comparative Experimental Study. In *SIGMETRICS*, 2014.
- [19] S. M. Khan et al. The Efficacy of Error Mitigation Techniques for DRAM Retention Failures: A Comparative Experimental Study. In *SIGMETRICS*, 2014.
- [20] Y. Kim et al. Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors. In *ISCA*, 2014.
- [21] Y. Kim et al. Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors. In *ISCA*, 2014.
- [22] E. Kultursay et al. Evaluating STT-RAM As An Energy-Efficient Main Memory Alternative. In *ISPASS*, 2013.
- [23] B. C. Lee et al. Architecting Phase Change Memory as a Scalable DRAM Alternative. In *ISCA*, 2009.
- [24] B. C. Lee et al. Phase Change Memory Architecture and The Quest for Scalability. *CACM*, 2010.
- [25] B. C. Lee et al. Phase Change Technology and the Future of Main Memory. *IEEE Micro*, 2010.
- [26] J. Liu et al. RAIDR: Retention-Aware Intelligent DRAM Refresh. In *ISCA*, 2012.
- [27] J. Liu et al. An Experimental Study of Data Retention Behavior in Modern DRAM Devices: Implications for Retention Time Profiling Mechanisms. In *ISCA*, 2013.
- [28] J. Liu et al. An Experimental Study of Data Retention Behavior in Modern DRAM Devices: Implications for Retention Time Profiling Mechanisms. In *ISCA*, 2013.
- [29] Y. Low et al. Distributed GraphLab: A Framework for Machine Learning and Data Mining in the Cloud. *PVLDB*, 2012.
- [30] Y. Luo et al. Characterizing Application Memory Error Vulnerability to Optimize Datacenter Cost via Heterogeneous-Reliability Memory. In *DSN*, 2014.
- [31] J. Meza et al. A Case for Efficient Hardware/Software Cooperative Management of Storage and Memory. In *WEED*, 2013.
- [32] J. Meza et al. A Large Scale Study of Flash Errors in the Field. In *SIGMETRICS*, 2015.
- [33] J. Meza et al. Revisiting Memory Errors in Large-Scale Production Data Centers: Analysis and Modeling of New Trends from the Field. In *DSN*, 2015.
- [34] O. Mutlu. Memory Scaling: A Systems Architecture Perspective. In *MEMCON*, 2013.
- [35] O. Mutlu. Memory Scaling: A Systems Architecture Perspective. In *IMW*, 2013.
- [36] M. Qureshi et al. AVATAR: A Variable-Retention-Time (VRT) Aware Refresh for DRAM Systems. In *DSN*, 2015.
- [37] M. K. Qureshi et al. Scalable High Performance Main Memory System Using Phase-Change Memory Technology. In *ISCA*, 2009.
- [38] M. K. Qureshi et al. Scalable High Performance Main Memory System Using Phase-Change Memory Technology. In *ISCA*, 2009.
- [39] V. J. Reddi et al. Web Search Using Mobile Cores: Quantifying and Mitigating the Price of Efficiency. In *ISCA*, 2010.
- [40] B. Schroeder et al. DRAM Errors in the Wild: A Large-Scale Field Study. In *SIGMETRICS Performance*, 2009.
- [41] V. Sridharan et al. Feng Shui of Supercomputer Memory: Positional Effects in DRAM and SRAM Faults. In *SC*, 2013.