

A Realistic Evaluation of Memory Hardware Errors and Software System Susceptibility

Xin Li Michael C. Huang Kai Shen Lingkun Chu
University of Rochester Ask.com
{xinli@ece, mihuang@ece, kshen@cs}.rochester.edu lchu@ask.com

Technical Report #949
Department of Computer Science, University of Rochester

September 3, 2009

Abstract

Memory hardware reliability is an indispensable part of whole-system dependability. Its importance is evidenced by a plethora of prior research work studying the impact of memory errors on software systems. However, the absence of solid understanding of the error characteristics prevents software system researchers from making well reasoned assumptions, and it also hinders the careful evaluations over different choices of fault tolerance design. In this paper, we present our realistic memory hardware error traces collected from production computer systems with more than 800 GB memory for around nine months. Based on the traces (including detailed information on the error addresses and patterns), we explore the implications of different hardware ECC protection schemes so as to identify the most common error causes and approximate error rates exposed to the software level. Lastly, we investigate the software system susceptibility to some major error causes, with the particular goal to validate, question, and augment results of prior system studies.

Key Words: reliability, memory hardware errors, software/hardware interaction.

1 Introduction

Memory hardware errors are an important threat to computer system reliability [31] as VLSI technologies continue to scale [6]. Past case studies suggested that hardware faults [29, 34] and particularly memory errors [24, 32] are significant contributing factors to whole-system failures. Understanding memory hardware errors is an important component in developing an overall system dependability strategy. Recent software system studies have attempted

to examine the impact of memory hardware errors on computer system reliability [11, 23] and security [13]. Software system countermeasures to these errors have also been investigated [27].

Despite its importance, our collective understanding about the rate, pattern, impact, and scaling trends of memory hardware errors is still somewhat fragmented and incomplete. For instance, while transient (or soft) errors received a thorough treatment in the literature [3, 26, 35–37], non-transient errors (including permanent and intermittent errors) have received less attention. Existing research on non-transient hardware failures often utilize synthetic error models [21]. The scarcity of realistic error traces is partly due to the fact that collecting field data requires access to large-scale facilities and these errors do not lend themselves to accelerated tests as transient errors do [37]. The two studies of non-transient errors that we are aware of [10, 30] provide no result on specific error locations and patterns.

In an effort to acquire valuable error statistics in real-world environments, we have monitor memory hardware errors in three groups of computers—specifically, a rack-mounted Internet server farm with more than 200 machines, about 20 university desktops, and 70 PlanetLab machines. We have collected error tracking results on over 800 GB memory for around nine months (from November 30, 2006 to September 11, 2007). Our error traces, probably the only public memory hardware error traces with detailed error addresses and patterns, are available through a link at the USENIX computer failure data repository [9].

One important discovery from our error traces is that non-transient errors are at least as significant a source of reliability concern as transient errors. In theory, permanent hardware errors, whose symptoms persist over time, are easier to detect. Consequently they ought to present only a minimum threat to system reliability in an ideally-

maintained environment. However, some non-transient errors are intermittent [10] (*i.e.*, whose symptoms are unstable at times) and they are not necessarily easy to detect. Further, the system maintenance is hardly perfect, particularly for hardware errors that do not trigger obvious system failures. Given our discovery of non-transient errors in real-world production systems, a holistic dependability strategy needs to take into account their presence and error characteristics.

We conduct trace-driven studies to understand hardware error manifestations and their impact on the software system. First, we extrapolate the collected traces into general statistical error manifestation patterns. We then perform Monte Carlo simulations to learn the error rate and particularly error causes under different memory protection mechanisms (*e.g.*, single-error-correcting ECC or stronger Chipkill ECC [12]). To achieve high confidence, we also study the sensitivity of our results to key parameters of our simulation model.

Further, we use a virtual machine-based error injection approach to study the error susceptibility of real software systems and applications. In particular, we discovered the previous conclusion that most memory hardware errors do not lead to incorrect software execution [11, 23] is not appropriate for non-transient memory errors. We also validated the failure oblivious computing model [28] using our web server workload with injected non-transient errors.

The rest of this paper is organized as follows. Section 2 provides a brief background on memory errors. Section 3 presents our collected raw error data from production systems. Section 4 analyzes error manifestation rate and leading error causes using Monte Carlo simulations. Section 5 then evaluates the software system susceptibility to memory hardware errors using error injection experiments. Finally, Section 6 concludes the paper.

2 Background

Terminology In general, a fault is the cause of an error, and errors lead to service failures [20]. Precisely defining these terms (“fault”, “error”, and “failure”), however, can be “surprisingly difficult” [2], as it depends on the notion of the system and its boundaries. For instance, the consequence of reading from a defective memory cell (obtaining an erroneous result) can be considered as a *failure* of the memory subsystem, an *error* in the broader computer system, or it may not lead to any failure of the computer system at all if it is masked by subsequent processing. In our discussion, we use error to refer to the incidence of having incorrect memory content. The root cause of an error is the fault, which can be a particle impact, or defects in some part of the memory circuit. Note that an error does not manifest (*i.e.*, it is a *latent error*) until the corrupt location is accessed.

An error may involve more than a single bit. Specifically, we count all incorrect bits due to the same root cause as part of one error. This is different from the concept of a multi-bit error in the ECC context, in which case the multiple incorrect bits must fall into a single ECC word. To avoid confusions we call these errors word-wise multi-bit instead.

Transient memory errors are those that do not persist and are correctable by software overwrites or hardware scrubbing. They are usually caused by temporary environmental factors such as particle strikes from radioactive decay and cosmic ray-induced neutrons. *Non-transient* errors, on the other hand, are often caused (at least partially) by inherent manufacturing defect, insufficient burn-in, or device aging [6]. Once they manifest, they tend to cause more predictable errors as the deterioration is often irreversible. However, before transitioning into permanent errors, they may put the device into a marginal state causing apparently *intermittent* errors.

Memory ECC Computer memories are often protected by some form of *parity-check code*. In a parity-check code, information symbols within a word are processed to generate *check* symbols. Together, they form the coded word. These codes are generally referred to as ECC (error correcting code). Commonly used ECC codes include SECDED and chipkill.

SECDED stands for *single-error correction, double-error detection*. Single error correction requires the code to have a Hamming distance of at least 3. In binary codes, it can be easily shown that r bits are needed for $2^r - 1$ information bits. For double-error detection, one more check bit is needed to increase the minimum distance to 4. The common practice is to use 8 check bits for 64 information bits forming a 72-bit ECC word as these widths are used in current DRAM standards (*e.g.*, DDR2).

Chipkill ECC is designed to tolerate word-wise multi-bit errors such as those caused when an entire memory device fails [12]. Physical constraints dictate that most memory modules have to use devices each providing 8 or 4 bits to fill the bus. This means that a chip-fail tolerant ECC code needs to correct 4 or 8 adjacent bits. While correcting multi-bit errors in a word is theoretically rather straightforward, in practice, given the DRAM bus standard, it is most convenient to limit the ECC word to 72 bits, and the 8-bit parity is insufficient to correct even a 4-bit symbol. To address this issue, one practice is to reduce the problem to that of single-bit correction by spreading the output of, say, 4 bits to 4 independent ECC words. The trade-off is that a DIMM now only provides 1/4 of the bits needed to fill the standard 64-data-bit DRAM bus, and thus a system needs a minimum of 4 DIMMs to function. Another approach is to use *b-adjacent* codes with much more involved matrices for parity generation and checking [7]. Even in this case, a typ-

ical implementation requires a minimum of 2 DIMMs. Due to these practical issues, chipkill ECC remains a technique used primarily in the server domain.

3 Realistic Memory Error Collection

Measurement results on memory hardware errors, particularly transient errors, are available in the literature. Ziegler *et al.* from IBM suggested that cosmic rays may cause transient memory bit flips [35] and did a series of measurements from 1978 to 1994 [26, 36, 37]. In a 1992 test for a vendor 4Mbit DRAM, they reported the rate of 5950 failures per billion device-hour. Published results on non-transient memory errors are few [10, 30] and they provide little detail on error addresses and patterns, which are essential for our analysis.

To enable our analysis on error manifestation and software susceptibility, we make efforts to collect realistic raw error rate and patterns on today’s systems. Our efforts are in two fronts. First, we perform long-term monitoring on large, non-biased sets of production computer systems. Second, we pursue outside reports of potential memory failure manifestation and investigate the error natures and patterns. Due to the rareness of memory hardware errors, the error collection can require enormous efforts. A general understanding of memory hardware errors is likely to require the collective and sustained effort from the research community as a whole. We are not attempting such an ambitious goal in this study. Instead, our emphasis is on the *realism* of our production system error collection. As such, we do not claim general applicability of our results.

3.1 Production System Error Monitoring

We monitor memory errors in three environments—a set of 212 production machines in a server farm at Ask.com [1], about 20 desktop computers at Rochester computer science department, and around 70 wide-area-distributed Planet-Lab machines. Preliminary monitoring results (of shorter monitoring duration, focusing exclusively on transient errors, with little result analysis) were reported in another paper [22]. Here we provide an overview of our latest monitoring results on all error types. Due to factors such as machine configuration, our access privileges, and load, we obtained uneven amount of information from the three error monitoring environments. Most of our results were acquired from the large set of server farm machines, where we have access to the memory chipset’s internal registers and can monitor the ECC-protected DRAM of all machines continuously. Below we focus our result reporting on the data obtained in this environment.

All 212 machines from the server farm use Intel E7520 chipset as memory controller hub [18]. Most machines have

4 GB DDR2 SDRAM. Intel E7520 memory controller is capable of both SECDED or Chipkill ECC. In addition to error detection and correction, the memory controller attempts to log some information about memory errors encountered. Unfortunately, this logging capability is somewhat limited—there are only two registers to track the addresses of two distinct errors. These registers will only capture the first two memory errors encountered. Any subsequent errors will not be logged until the registers are reset. Therefore, we periodically (once per hour) probe the memory controller to read out the information and reset the registers. This probing is realized through enhancements of the memory controller driver [5], which typically requires the administrative privilege on target machines.

Recall that when a memory cell’s content is corrupted (creating a latent error), the error will not manifest to our monitoring system until the location is accessed. To help expose these latent errors, we enable hardware memory scrubbing—a background process that scans all memory addresses to detect and correct errors. The intention is to prevent errors from accumulating into more severe forms (*e.g.*, multi-bit) that are no longer correctable. It is typically performed at a low frequency (*e.g.*, 1.5 hours for every 1 GB) [18] to minimize the energy consumption and contention with running applications. Note that scrubbing does not help expose *faults*—writing varying values into memory does that. Since we monitored the machines for an extended period of time (9 months), the natural usage of the machines is likely to have exposed most (if not all) faults.

We collected error logs for a period of approximately 9 months (from November 30, 2006 to September 11, 2007). In the first 2 months we observed errors on 11 machines. No new errors were seen for 6 months and then 1 more erroneous machine appeared in the most recent month of our monitoring. We choose 6 erroneous machines with distinct error patterns and Figure 1 demonstrates how the errors are laid out on the physical memory arrays. Based on observed patterns, all four memory error modes (single-cell, row, column, and whole-chip [4]) appear in our log. Specifically, M10 contains a single cell error, M7 and M12 represent a row error and a column error respectively. Finally, for machine M8, the errors are spread all over the chip which strongly suggests faults in the chip-wide circuitry rather than individual cells, rows, or columns. Based on the pattern of error addresses, we categorize all error instances into appropriate modes shown in Table 1.

While the error-correction logic can detect errors, it cannot tell whether an error is transient or not. We can, however, make the distinction by continued observation—repeated occurrences of error on the same address are virtually impossible to be external noise-induced transient errors as they should affect all elements with largely the same probability. We can also identify non-transient errors by

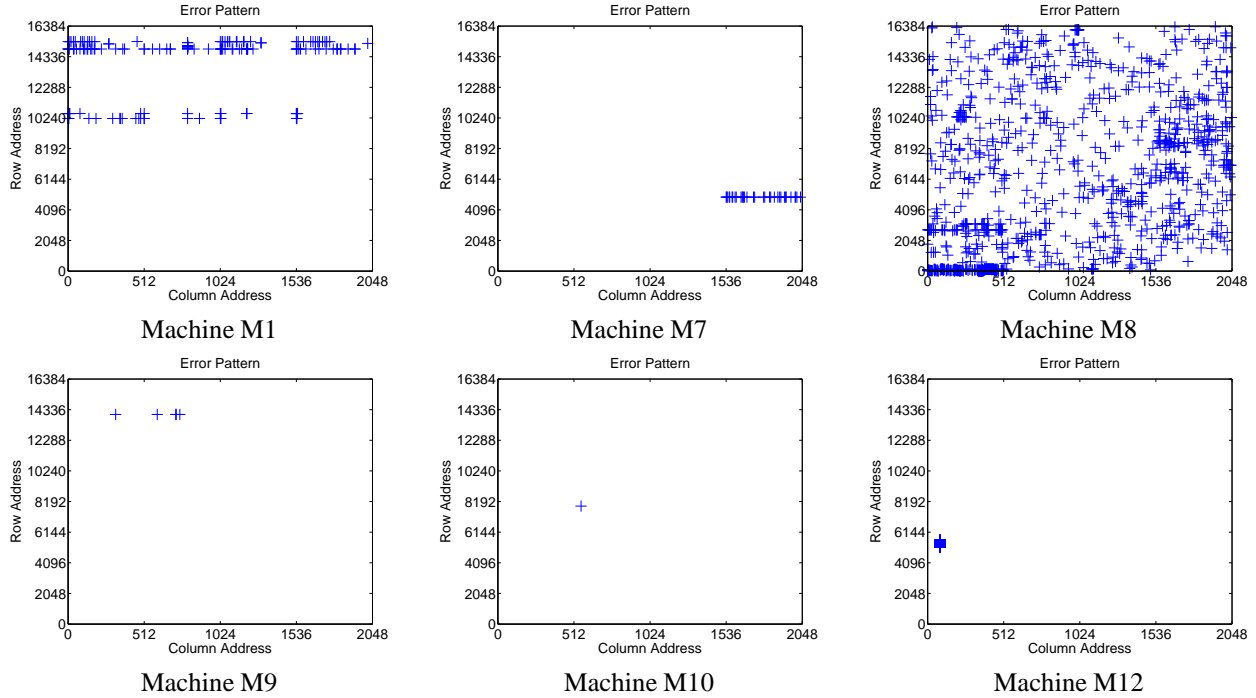


Figure 1. The visualization of error patterns on physical memory devices. Each cross represents an erroneous cell at its row/column addresses. The system address to row/column address translation is obtained from the official Intel document [18].

Machine	Single-cell	Row	Column	Whole-chip
M1	2	11	1	
M2		1		
M3	1 (transient)			
M4	1			
M5	1 (transient)			
M6				1
M7		1		
M8				1
M9		1		
M10	1			
M11	1			
M12			1	
Total	7 (2 transient)	14	2	2

Table 1. Collected errors and their modes.

recognizing known error modes related to inherent hardware defects: single-cell, row, column, and whole-chip [4]. For instance, memory row errors will manifest as a series of errors with addresses on the same row. Some addresses on this row may be caught on the log only once. Yet, the cause of that error is most likely non-transient if other cells on the same row indicate non-transient errors (logged multiple times). Take M9 in Figure 1 as an example. Altogether, there are five distinct error addresses recorded in our trace, two of which showed up only once and the rest were all recorded multiple times. Since they happen on the same row, it is highly probable that they are all due to defects in

places like the word line. We count them as a row error.

3.2 Investigation of Reported Failures

We also investigate outside reports of potential memory failures. Because these reports come from biased machine samples (those with suspicious failure symptoms), such error results are not amenable for general statistical analysis. Rather, it only serves the purpose to illustrate the broad scope of systems on which memory errors may manifest.

Here we provide an example of such investigation. We followed a local student’s report of memory failure on a medical System-on-Chip platform. The faulty chip is a Microchip PIC18F452, equipped with three kinds of memory—32 KB program memory, 256 bytes EEPROM for static data, and 1.5 KB SRAM for volatile data. The chip was used to monitor heart rate of neonates and it reported mysterious rate drops of 64. Using the “in-circuit debugger”, we were able to attribute the failure cause to a memory bit stuck at ‘1’ at the 23rd byte in the SRAM.

4 Error Manifestation Analysis

We analyze how device-level errors would be exposed to software. We are interested in the error manifestation rates and patterns (*e.g.*, multi-bits or single-bit) as well as leading causes for manifested errors. We explore results

DRAM technology	DDR2
DIMM No. per machine	4
Device No. per DIMM	18
Device data width	x4
Row/Column/Bank No.	$2^{14}/2^{11}/4$
Device capacity	512 Mb
Capacity per machine	4 GB
ECC capability	None, SECEDED, Chipkill

Table 2. Memory configuration for our server farm machines.

under different memory protection schemes. This is useful since Chipkill ECC represents a somewhat extreme trade-off between reliability and other factors (*e.g.*, performance and energy consumption) and may remain a limited-scope solution. In our memory chipset (Intel E7520) for example, to provide the necessary word length, the Chipkill design requires two memory channels to operate in a lock-stepping fashion, sacrificing throughput and power efficiency.

4.1 Evaluation Methodology

We use a discrete-event simulator to conduct Monte-Carlo simulations to derive properties of manifested errors. We simulate 500 machines with the exact configuration as the Ask.com servers in Section 3. The detailed configuration is shown in Table 2. We first use the error properties extracted from our data to generate error instances in different memory locations in the simulated machines. Then we simulate different ECC algorithms to obtain a trace of manifested memory errors as the output. Our analysis here does not consider software susceptibility to manifested errors, which will be examined in Section 5. Below, we describe several important aspects of our simulation model, including temporal error distributions, device-level error patterns, and the repair maintenance model.

- We consider transient and non-transient errors separately in terms of temporal error distribution. Since transient errors are mostly induced by random external events, it is well established that their occurrences follow a memoryless exponential distribution. The cumulative distribution function of exponential distribution is $F(t) = 1 - e^{-\lambda t}$, which represents the probability that an error has already occurred by time t . The instantaneous error rate for exponential distribution is constant over time, and does not depend on how long the chip has been operating properly.

The non-transient error rate follows a “bathtub” curve with a high, but declining rate in the early “infant mortality” period, followed by a long and stable period with a low rate, before rising again when device wear-out starts to take over. Some study has also suggested that improved manufacturing techniques com-

binated with faster upgrade of hardware have effectively made the wear-out region of the curve irrelevant [25]. In our analysis, we model 16 months of operation and ignore aging or wear-out. Under these assumptions, we use the oft-used Weibull distributions which has the following cumulative distribution function: $F(t) = 1 - e^{-(t/\beta)^\alpha}$. The *shape parameter* α controls how steep the rate decreases, and the *scale parameter* β determines how “stretched out” the curve is. Without considering the wear-out region, the shape parameter in the Weibull distribution is no more than 1.0, at which point the distribution degenerates into an exponential distribution. The temporal error occurrence information in our data suggested a shape parameter of 0.11.

- We then consider device-level error patterns. For transient errors, prior studies and our own observation all point to the single-bit pattern. For non-transient errors, we have the 10 distinct patterns in our trace as templates. When a non-transient error is to be generated, we choose one out of these templates in a uniformly random fashion. There is a problem associated with using the exact template patterns—error instances generated from the same templates are always injected on the same memory location and thus they would always be aligned together to cause an uncorrectable error in the presence of ECC. To address this problem, we shift the error location by a random offset each time we inject an error instance.
- Our model requires a faulty device repair maintenance strategy. We employ an idealized “reactive” repair without preventive maintenance. We assume an error is detected as soon as it is exposed to the software level. If the error is diagnosed to be non-transient, the faulty memory module is replaced. Otherwise the machine will undergo a reboot. In our exploration, we have tried two other maintenance models that are more proactive. In the first case, hardware scrubbing is turned on so that transient errors are automatically corrected. In the second case, we further assume that the memory controller notifies the user upon detecting a correctable non-transient error so that faulty memory modules can be replaced as early as possible. We found these preventive measures have a negligible impact on our results. We will not consider these cases in this paper.

Below, we provide evaluation results using the above described model (Section 4.2). Due to the small number of errors in the collected error trace, the derived rate and temporal manifestation pattern may not provide high statistical confidence. To achieve high confidence, we further study the sensitivity of our results to two model parameters—the

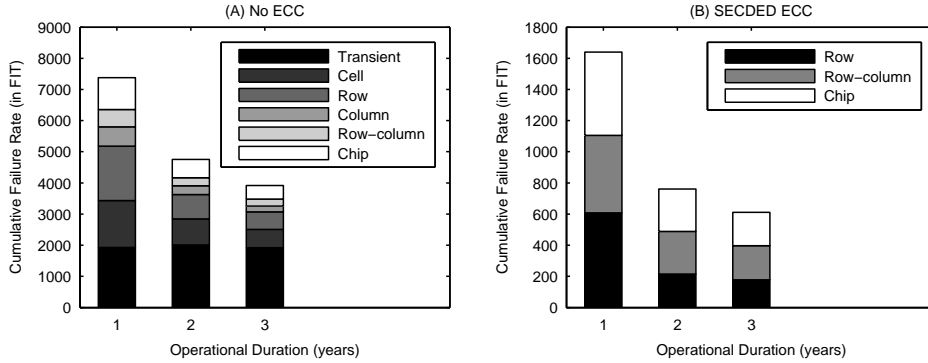


Figure 2. Failure rates and breakdown causes for no ECC and SECCDED ECC, with varying machine operational durations.

Weibull distribution shape parameter (Section 4.3) and the error rate over time (Section 4.4).

4.2 Base Results

Here we present the simulation results on failures. The failure rates are computed as the average of the simulated operational duration. We describe our results under different memory protection schemes.

Figure 2(A) illustrates the failure rates and the breakdown of the causes when there is no ECC protection. In this case, any error will be directly exposed to software and cause a failure. As a result, we can study the errors in isolation. With our measurement, the transient error rate is 2006 FIT¹ for each machine’s memory system. Depending on the operational time of the machines, the average non-transient error rates would vary, and so are the corresponding failure rates. Overall, for machines without ECC support, both transient and non-transient errors contribute to the overall error rate considerably.

SECCDED ECC can correct any word-wise single-bit errors. Of the errors in our trace, this capability will correct all but one whole-chip error, one row error, and one row-column error. These three cases all have multiple erroneous bits (due to the same root cause) in one ECC word, preventing ECC correction. Theoretically, a failure can also occur when multiple independent single-bit errors happen to affect the same ECC word (such as when a transient error occurs to an ECC word already having a single-bit non-transient error). However, since errors are rare in general, such combination errors are even less probable. In our simulations, no such instance has been encountered. Figure 2(B) summarizes the simulation results.

¹FIT is a commonly used unit to measure failure rates and 1 FIT equals one failure per billion device-hour. To put the numbers into perspectives, IBM’s target FIT rates for servers are 114 for undetected (or silent) data corruption, 4500 for detected errors causing system termination, and 11400 for detected errors causing application termination [8]. Note that these rates are for the whole system including all components.

When using the Chipkill ECC, as expected, the memory system becomes very resilient. We did not see any uncorrected errors. This result echoes the conclusion of [12].

4.3 Shape Parameter Sensitivity

To reach high confidence in our results, we consider a wide range of the Weibull shape parameters for the non-transient error temporal distribution and study the sensitivity of our results to this parameter. We use a machine operational duration of 16 months, which is the age of the Ask.com servers at the end of our data collection.

Prior failure mode studies in computer systems [14, 34], spacecraft electronics [15], electron tubes [19], and integrated circuits [17] pointed to a range of shape parameter values in 0.28–0.80. Given this and the fact that the Weibull distribution with shape parameter 1.0 degenerates to an exponential distribution, we consider the shape parameter range of 0.1–1.0 in this sensitivity study.

In both ECC mechanisms, the non-transient error rate depends on the Weibull shape parameter. The lower the shape parameter, the faster the error rate drops, and the lower the total error rate for the entire period observed. Note that the transient error rate also fluctuates a little because of the non-deterministic nature of our Monte-Carlo simulation. It’s apparent the change of transient error rates do not correlate with the shape parameter. For no-ECC, as Figure 3(A) shows, for machines in their first 16 months of operation, the difference caused by the wide ranging shape parameter is rather insignificant.

In the case of SECCDED shown in Figure 3(B), the impact of the Weibull shape parameter is a bit more pronounced than in the case of no ECC but is still relatively insignificant. Also, even though error rates are significantly reduced by SECCDED, they are still within a factor of about five from those without ECC.

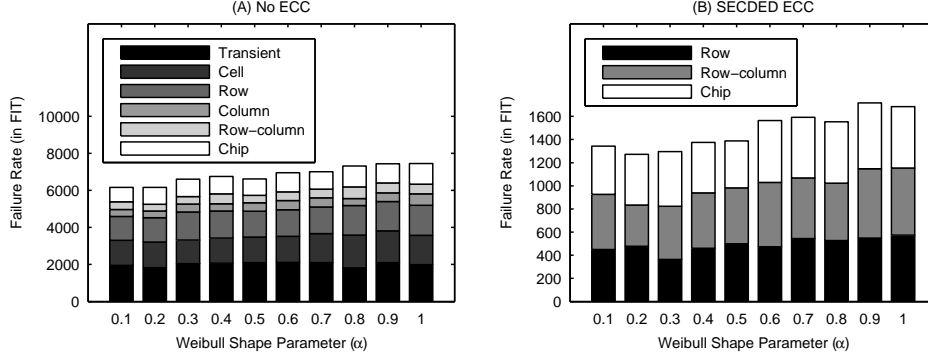


Figure 3. Failure rates and breakdown causes for no ECC and SECDED ECC, with varying Weibull shape parameter.

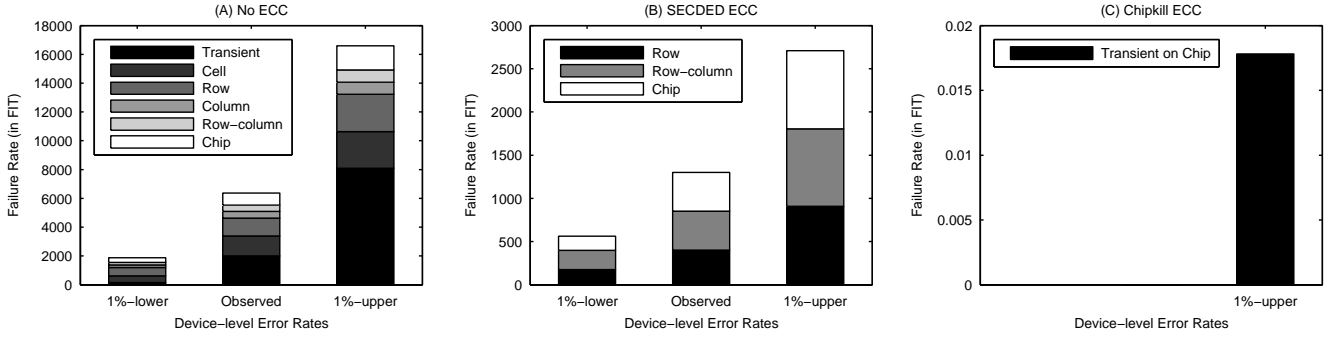


Figure 4. Manifested errors when input device-level error rates are the originally observed and 1%-lower/upper-bounds.

4.4 Statistical Error Rate Bounds

Due to the small number of device-level errors in our trace, the observed error rate may be quite different from the intrinsic error rate of our monitored system. To account for such inaccuracy, we use the concept of *p-value bounds* to provide a range of possible intrinsic error rates with statistical confidence.

For a given probability p , the p -value upper bound (λ_u) is defined as the intrinsic error rate under which $Pr\{X \leq n\} = p$. Here n is the actual number of errors observed in our experiment. X is the random variable for the number of errors occurring in an arbitrary experiment of the same time duration. And likewise, the p -value lower bound (λ_l) is the intrinsic error rate under which $Pr\{X \geq n\} = p$. A very small p indicates that given n observed errors, it is improbable for the actual intrinsic error rate λ to be greater than λ_u or less than λ_l .

Given p , the probability distribution of random variable X is required to calculate the p -value for our data. Thankfully, when the memory chips are considered identical, we can avoid this requirement. This is because in any time interval, their probability of having an error is the same, say q . Let N be the total number of memory chips operating, then the actual number of errors happening in this period, X , will be a random variable which conforms to binomial

distribution: $P_{N,q}\{X = k\} = \binom{N}{k} q^k (1 - q)^{N-k}$. When N is very large (we simulated thousands of chips), we can approximate by assuming N approaches infinity. In this case the binomial distribution will turn into Poisson distribution. For the ease of calculation, we shall use the form of Poisson distribution: $P_\lambda\{X = k\} = \frac{e^{-\lambda} \lambda^k}{k!}$, where $\lambda = q \cdot N$ is the expectation of X .

Based on the analysis above and the observed error rates, we have calculated the 1% upper and lower bounds. For instance, the transient error rate in non-ECC memory system is 2006 FIT as mentioned earlier. The corresponding 1%-upper-bound and 1%-lower-bound are 8429 FIT and 149 FIT respectively. The bounds on the various manifested error rates, derived from different raw error rates, are shown in Figure 4. From left to right, the bars show the 1%-lower-bound, the originally observed rate, and the 1%-upper-bound. As can be seen, for manifestations caused by non-transient errors, the two 1% bounds are roughly 2x to either direction of the observed rate. This range is narrow enough that there is little impact to the qualitative conclusions.

For Chipkill ECC, the 1%-upper-bound offers a better chance to observe failures in the outcome of our simulation. With this increased rate, we finally produced a few failure instances (note there were none for Chipkill in the base sim-

ulations done in previous sub-sections). The patterns of the failures are shown in Figure 4(C). All of the failures here are caused by a transient error hitting an existing non-transient chip error.

4.5 Summary

We summarize our major findings of this part of the study: 1) In terms of the absolute failure rate, with no ECC protection, error rates are at the level of thousands of FIT per machine. SECDED ECC lowers the rates to the neighborhood of 1000 FIT per machine. Chipkill ECC renders failure rates virtually negligible. 2) Non-transient errors are significant (if not dominant) causes for all cases that we evaluated. Particularly on SECDED ECC machines, manifested failures tend to be caused by row errors, row-column errors, and whole-chip errors. 3) Word-wise multi-bit failures are quite common.

5 Software System Susceptibility

A memory error that escaped hardware ECC correction is exposed to the software level. However, its corrupted memory value may or may not be consumed by software programs. Even if it is consumed, the software system and applications may continue to behave correctly if such correctness does not depend on the consumed value. Now we shift our attention to the susceptibility of software systems and applications to memory errors. Specifically, we inject the realistic error patterns from our collected traces and observe the software behaviors. Guided by the conclusion of Section 4, we also take into account the shielding effect of ECC algorithms.

There is a rich body of prior research on software system reliability or security regarding memory hardware errors [11, 13, 21, 23, 27, 28]. One key difference between these studies and ours is that all of our analysis and discussions ultimately root in the realism of our collected error trace. In this section, we tailor our software susceptibility evaluation in the context of recent relevant research with the hope of validating, questioning, or augmenting prior results.

5.1 Methodology of Empirical Evaluation

To run real software systems on injected error patterns, we must accomplish the following goals. First, every read access to a faulty location must be supplied with an erroneous value following the injection pattern. This can be achieved by writing the erroneous value to each individual faulty address at the time of injection. Second, for every write access to a faulty location, if the error is non-transient, we must guarantee the erroneous value is restored right after the write. The injection is then followed by error man-

ifestation bookkeeping. The bookkeeping facility has to be informed whenever a faulty address is accessed so that it would log some necessary information. The key challenge of such error injection and information logging is to effectively track and manipulate all the accesses to locations injected with errors (or *tracked locations*).

We base our tracking method on a *Page Access Control* approach proposed by [33], with some performance enhancements. The basic idea of this approach is to identify the pages that contain the tracked memory locations, and then apply page access protection in the page table. Once the control is applied, an access to a tracked location will raise a page protection fault so that the system will get notified of the access. To make the access complete and resume the application, we simply remove the page access control temporarily and single-step the current instruction. At the end of the single-stepping, we restore the control.

If the error injection and monitoring mechanisms are built into the target system itself (as in [23]), these mechanisms may not behave reliably in the presence of injected memory errors. To avoid this potential problem, we utilize a virtual machine-based architecture in which the target system runs within a hosted virtual machine while the error injection and monitoring mechanisms are built in the underlying virtual machine monitor. We enable the shadow page table mode in the virtual machine memory management. Error injections only affect the shadow page tables while page tables within the target virtual machine are not affected. In order to understand software system susceptibility to memory hardware errors, we log certain information every time an error is activated. Specifically, we record the access type (read or write), access mode (kernel or user), and the program counter value. For kernel mode accesses, we are able to locate specific operating system functions from the program counter values.

Our experimental environment employs Xen 3.0.2 and runs the target system in a virtual machine with Linux 2.6.16 operating system. We examine three applications in our test: 1) the Apache web server running the static request portion of the SPECweb99 benchmark with around 2 GB web documents. 2) MCF from SPEC CPU2000—a memory-intensive vehicle scheduling program for mass transportation; and 3) compilation and linking of the Linux 2.6.23 kernel. The first is a typical server workload while the other two are representative workstation workloads (in which MCF is CPU-intensive while kernel build involves significant I/O).

5.2 Evaluation on Failure Severity

Two previous studies [11, 23] investigated the susceptibility of software systems to transient memory errors. They reached similar conclusions that overall, memory errors do

not pose a significant threat to software systems. In particular, Messer *et al.* [23] discovered that of all the errors they injected, on average 20% were accessed, among which 74% are overwritten before being really consumed by the software. In other words, only 5% of the errors would cause abnormal software behaviors. However, these studies limited their scope for single-bit transient errors only. Our findings in Section 4 show non-transient errors are also a significant cause of memory failures. When these errors are taken into account, the previous conclusions may not stand intact. For example, non-transient errors may not be overwritten, and as a result, a portion of the 74% overwritten errors in [23] would have been consumed by the software system if they had been non-transient.

Table 3 summarizes the execution results of our three benchmark applications when non-transient errors are injected. Since our applications all finish in a short time (a few minutes), we consider these non-transient errors as permanent during the execution. In total we had 12 different error patterns. M3 and M5 are transient errors and therefore we do not include them in this result. M8 is so massive that as soon as it is injected, the OS crashes right away. We also exclude it from our results.

The table includes results for both cases of no ECC and SECDED ECC. Since errors are extremely rare on Chipkill machines (see conclusions of Section 4), here we do not provide results for Chipkill. For no ECC, briefly speaking, out of the 27 runs, 13 have accessed memory errors and 8 did not finish with expected correct results. This translates to 48% of the errors are activated and 62% of the activated errors do lead to incorrect execution of software systems. In the SECDED case, single-bit errors would be corrected. Most errors (except M1 and M7) are completely shielded by the SECDED ECC. However, for the six runs with error patterns M1 and M7, five accessed the errors and subsequently caused abnormal behaviors.

Overall, compared to results in [23], non-transient errors evidently do cause more severe consequences to software executions. The reason for the difference is twofold— 1) non-transient errors are not correctable by overwriting and 2) unlike transient errors, non-transients sometimes involve a large number of erroneous bits. To demonstrate reason #1, we show in Table 4, when these errors are turned into transient ones (meaning they can be corrected by overwritten values), quite a few of the execution runs would finish unaffected.

5.3 Validation of Failure-Oblivious Computing

In this section we attempt to validate the concept of failure-oblivious computing [28] with respect to memory hardware errors. It is based on the premise that in server

Application	Kernel build	Web server	MCF
No ECC			
M1 (Row-Col error)	AC	WO	AC
M2 (Row error)		OK	
M4 (Bit error)		OK	
M6 (Chip error)	AC	KC	WO
M7 (Row error)		WO	WO
M9 (Row error)		OK	
M10 (Bit error)		OK	
M11 (Bit error)			
M12 (Col error)		WO	
SECDED ECC			
M1 (Row-Col error)	AC	WO	WO
M7 (Row error)		WO	WO

Table 3. Error manifestation for each of our three applications. The abbreviations in the table should be interpreted as follows, with descending manifestation intensity: KC—kernel crash; AC—application crash; WO—wrong output; OK—program runs correctly. The blank cells indicate the error was not accessed at all.

Application	Kernel build	Web server	MCF
No ECC			
M1 (Row-Col error)	OK	WO	AC
M2 (Row error)		OK	
M4 (Bit error)		OK	
M6 (Chip error)	OK	KC	OK
M7 (Row error)		WO	OK
M9 (Row error)		OK	
M10 (Bit error)		OK	
M11 (Bit error)			
M12 (Col error)		WO	
SECDED ECC			
M1 (Row-Col error)	OK	WO	OK
M7 (Row error)		WO	OK

Table 4. Error manifestation for each of our three applications, when the errors are made transient (thus correctable by overwrites). Compared to Table 3, many of the runs are less sensitive to transient errors and exhibit no mis-behavior at the application level.

workloads, error propagation distance is usually very small. When memory errors occur (mostly they were referring to out-of-bound memory accesses), a failure-oblivious computing model would discard the writes and supply the read with arbitrary values and try to proceed. In this way the error occurred will be confined within the local scope of a request and the server computation can be resumed without being greatly affected.

The failure-oblivious concept may also apply to memory hardware errors. It is important to know what the current operating system does in response to memory errors. Without ECC, the system is obviously unaware of any memory errors going on. Therefore it is truly failure-oblivious. With ECC, the system could detect some of the uncorrectable errors. At this point the system can choose to stop, or to con-

Application	No ECC	SECDED ECC
M1 (Row-Col error)	15	8
M2 (Row error)	0	0
M3 (Transient error)	0	0
M4 (Bit error)	0	0
M5 (Transient error)	0	0
M7 (Row error)	2	1
M9 (Row error)	0	0
M10 (Bit error)	0	0
M11 (Bit error)	0	0
M12 (Col error)	1	0

Table 5. Number of requests affected by the errors in SPECweb99-driven Apache web server. We request 14400 files in each run.

tinue execution (probably with some form of error logging). The specific choices are machine dependent.

For our web server workload, we check the integrity of web request returns in the presence of memory errors. Table 5 lists the number of requests with wrong contents for each non-transient error. The worst case is M1, which caused 15 erroneous request returns (or files with incorrect content). However, this is still a small portion (about 0.1%) in the total 14400 files we have requested. Our result suggests that, in our tested web server workload, memory-hardware-error-induced failures tend not to propagate very far. This shows the promise of applying failure-oblivious computing in the management of memory hardware errors for server systems.

5.4 Discussion on Additional Cases

Though error testing data from the industry are seldom published, modern commercial operating systems do advocate their countermeasures for faulty memory. Both IBM AIX [16] and Sun Solaris [32] have the ability to retire faulty memory when the ECC reports excessive correctable memory errors. Our results suggest that with ECC protection, the chances of errors aligning together to form an uncorrectable one is really low. However, this countermeasure could be effective against those errors that gradually develop into uncorrectable ones by themselves. Since our data does not have timestamps for most of the error instances, it is hard to verify how frequently these errors occur. On Chipkill machines [16], however, this countermeasure seems to be unnecessary since our data shows that without any replacement policy, Chipkill will maintain the memory failure rate at an extremely low level.

A previous security study [13] devised a clever attack that exploits memory errors to compromise the Java virtual machine (JVM). They fill the memory with pointers to an object of a particular class, and through an accidental bit flip, they hope one of the pointers can point to an object of another class. Obtaining a class A pointer actually pointing

to a class B object is enough to compromise the whole JVM. In particular, they also provided an analysis of the effectiveness of exploiting multi-bit errors [13]. It appears that they can only exploit bit flips in a region within a pointer word (in their case, bit 2:27 for a 32-bit pointer). In order for an error to be exploitable, all the bits involved must be in the region. The probability that they can exploit the error decreases with the number of erroneous bits in the word. Considering that the multi-bit errors in our collected error trace are mostly consecutive rather than distributed randomly, we can be quite optimistic about successful attacks.

Another previous study [27] proposed a method to protect critical data against illegal memory writes as well as memory hardware errors. The basic idea is that software systems can create multiple copies of their critical data. If a memory error corrupts one copy, a consistency check can detect and even correct such errors. The efficacy of such an approach requires that only one copy of the critical data may be corrupted at a time. Using our collected realistic memory error patterns, we can explore how the placement of multiple critical data copies affects the chance for simultaneous corruption. In particular, about half of our non-transient errors exhibit regular column or row-wise array patterns. Therefore, when choosing locations for multiple critical data copies, it is best to have them reside in places with different hardware row and column addresses (especially row addresses).

6 Conclusion

In this paper, we have presented a set of memory hardware error data collected from production computer systems with more than 800 GB memory for around 9 months. We discover a significant number of non-transient errors (typically in the patterns of row or column errors). Driven by the collected error patterns and taking into account various ECC protection schemes, we conducted a Monte Carlo simulation to analyze how errors manifest at the interface between the memory subsystem and software applications. Our basic conclusion is that non-transient errors comprise a significant portion of the overall errors visible to software systems. In particular, with the conventional ECC protection scheme of SECDED, transient errors will be almost eliminated while only non-transient memory errors may affect software systems and applications.

We also investigated the susceptibility of software system and applications to realistic memory hardware error patterns. In particular, we find that the earlier results that most memory hardware errors do not lead to incorrect software execution [11, 23] may not be valid, due to the unrealistic model of exclusive transient errors. At the same time, we provide a validation for the failure-oblivious computing model [28] on our web server workload with injected mem-

ory hardware errors.

Acknowledgments

We would like to thank Tao Yang and Alex Wong at Ask.com who helped us in acquiring administrative access to Ask.com Internet servers. We would also like to thank Howard David at Intel for kindly interpreting the memory error syndromes. This work was supported in part by the National Science Foundation (NSF) grants CCR-0306473, ITR/IIS-0312925, CNS-0509270, CNS-0615045, and CCF-0621472. Kai Shen was also supported by an NSF CAREER Award CCF-0448413 and an IBM Faculty Award.

References

- [1] Ask.com (formerly Ask Jeeves Search). <http://www.ask.com>.
- [2] A. Avizienis, J.-C. Laprie, B. Randell, and C. E. Landwehr. Basic concepts and taxonomy of dependable and secure computing. *IEEE Trans. Dependable Sec. Comput.*, 1(1):11–33, 2004.
- [3] R. Baumann. Soft errors in advanced computer systems. *IEEE Design and Test of Computers*, 22(3):258–266, May 2005.
- [4] M. Blaum, R. Goodman, and R. McEliece. The reliability of single-error protected computer memories. *IEEE Trans. on Computers*, 37(1):114–118, 1988.
- [5] EDAC project. <http://bluesmoke.sourceforge.net>.
- [6] S. Borkar. Designing reliable systems from unreliable components: the challenges of transistor variability and degradation. *IEEE Micro*, 25(6):10–16, Nov.–Dec. 2005.
- [7] D. Bossen. *b*-adjacent error correction. *IBM Journal of Research and Development*, 14(4):402–408, 1970.
- [8] D. Bossen. CMOS soft errors and server design. In *2002 Reliability Physics Tutorial Notes – Reliability Fundamentals*, pages 121.07.1 – 121.07.6, Dallas, Texas, Apr. 2002.
- [9] USENIX computer failure data repository. <http://cfdx.usenix.org>.
- [10] C. Constantinescu. Impact of deep submicron technology on dependability of VLSI circuits. In *Int'l Conf. on Dependable Systems and Networks*, pages 205–209, Bethesda, MD, June 2002.
- [11] C. da Lu and D. A. Reed. Assessing fault sensitivity in MPI applications. In *Supercomputing*, Pittsburgh, PA, Nov. 2004.
- [12] T. J. Dell. A white paper on the benefits of chipkill correct ECC for PC server main memory. *White paper*, 1997.
- [13] S. Govindavajhala and A. W. Appel. Using memory errors to attack a virtual machine. In *IEEE Symp. on Security and Privacy*, pages 154–165, Berkeley, CA, May 2003.
- [14] T. Heath, R. P. Martin, and T. D. Nguyen. Improving cluster availability using workstation validation. In *ACM SIGMETRICS*, pages 217–227, Marina del Rey, CA, June 2002.
- [15] H. Hecht and E. Fiorentino. Reliability assessment of spacecraft electronics. In *Annu. Reliability and Maintainability Symp.*, pages 341–346. IEEE, 1987.
- [16] D. Henderson, B. Warner, and J. Mitchell. IBM POWER6 processor-based systems: Designed for availability. *White paper*, 2007.
- [17] D. P. Holcomb and J. C. North. An infant mortality and long-term failure rate model for electronic equipment. *AT&T Technical Journal*, 64(1):15–31, January 1985.
- [18] Intel E7520 chipset datasheet: Memory controller hub (MCH). http://www.intel.com/design/chipsets/E7520_E7320/documentation.htm.
- [19] J. H. K. Kao. A graphical estimation of mixed Weibull parameters in life-testing of electron tubes. *Technometrics*, 1(4):389–407, Nov. 1959.
- [20] J. Laprie. Dependable computing and fault tolerance: Concepts and terminology. In *15th Int'l Symp. on Fault-Tolerant Computing*, pages 2–11, Ann Arbor, MI, June 1985.
- [21] M. Li, P. Ramachandran, S. K. Sahoo, S. V. Adve, V. S. Adve, and Y. Zhou. Understanding the propagation of hard errors to software and implications for resilient system design. In *13th Int'l Conf. on Architectural Support for Programming Languages and Operating Systems*, pages 265–276, Seattle, WA, Mar. 2008.
- [22] X. Li, K. Shen, M. Huang, and L. Chu. A memory soft error measurement on production systems. In *USENIX Annual Technical Conf.*, pages 275–280, Santa Clara, CA, June 2007.
- [23] A. Messer, P. Bernadat, G. Fu, D. Chen, Z. Dimitrijevic, D. J. F. Lie, D. Mannaru, A. Riska, and D. S. Milojevic. Susceptibility of commodity systems and software to memory soft errors. *IEEE Trans. on Computers*, 53(12):1557–1568, 2004.
- [24] B. Murphy. Automating software failure reporting. *ACM Queue*, 2(8):42–48, Nov. 2004.
- [25] F. R. Nash. *Estimating Device Reliability: Assessment of Credibility*. Springer, 1993. ISBN 079239304X.
- [26] T. J. O’Gorman, J. M. Ross, A. H. Taber, J. F. Ziegler, H. P. Muhlfield, C. J. Montrose, H. W. Curtis, and J. L. Walsh. Field testing for cosmic ray soft errors in semiconductor memories. *IBM J. of Research and Development*, 40(1):41–50, 1996.
- [27] K. Pattabiraman, V. Grover, and B. G. Zorn. Samurai: Protecting critical data in unsafe languages. In *Third EuroSys Conf.*, pages 219–232, Glasgow, Scotland, Apr. 2008.
- [28] M. Rinard, C. Cadar, D. Dumitran, D. M. Roy, T. Leu, and J. William S. Beebe. Enhancing server availability and security through failure-oblivious computing. In *6th USENIX Symp. on Operating Systems Design and Implementation*, pages 303–316, San Francisco, CA, Dec. 2004.
- [29] B. Schroeder and G. A. Gibson. A large-scale study of failures in high-performance computing systems. In *Int'l Conf. on Dependable Systems and Networks*, pages 249–258, Philadelphia, PA, June 2006.
- [30] B. Schroeder, E. Pinheiro, and W.-D. Weber. DRAM errors in the wild: a large-scale field study. In *ACM SIGMETRICS*, pages 193–204, 2009.
- [31] Sun Microsystems server memory failures. <http://www.forbes.com/global/2000/1113/0323026a.html>.

- [32] D. Tang, P. Carruthers, Z. Totari, and M. W. Shapiro. Assessment of the effect of memory page retirement on system RAS against hardware faults. In *Int'l Conf. on Dependable Systems and Networks*, pages 365–370, Philadelphia, PA, June 2006.
- [33] R. Wahbe. Efficient data breakpoints. In *5th Int'l Conf. on Architectural Support for Programming Languages and Operating Systems*, pages 200–212, Boston, MA, Oct. 1992.
- [34] J. Xu, Z. Kalbarczyk, and R. K. Iyer. Networked Windows NT system field failure data analysis. In *Pacific Rim Intl. Symp. on Dependable Computing*, pages 178–185, Hong Kong, China, Dec. 1999.
- [35] J. Ziegler and W. Lanford. Effect of cosmic rays on computer memories. *Science*, 206(16):776–788, Nov. 1979.
- [36] J. Ziegler, M. Nelson, J. Shell, R. Peterson, C. Gelderloos, H. Muhlfield, and C. Montrose. Cosmic ray soft error rates of 16-Mb DRAM memory chips. *IEEE Journal of Solid-State Circuits*, 33(2):246–252, Feb. 1998.
- [37] J. F. Ziegler et al. IBM experiments in soft fails in computer electronics (1978–1994). *IBM J. of Research and Development*, 40(1):3–18, 1996.