# Impact of Intermittent Faults on Nanocomputing Devices

Cristian Constantinescu
*Advanced Micro Devices Corp.*
*2950 E Harmony Road*
*CO 80528, USA*
*cristian.constantinescu@amd.com*

## Abstract

*Operation of semiconductor devices may be negatively affected by permanent, transient and intermittent faults. Nanocomputing devices are expected to experience higher error rates, in particular due to transient and intermittent faults. The errors induced by high energy particles, usually referred to as soft errors, have been extensively studied. However, similar errors may be induced by intermittent faults. This paper defines the permanent, intermittent and transient fault classes, emphasizing several malfunction causes, from manufacturing residues to ultra-thin oxide breakdown and timing violations. Error signatures, specific to intermittent faults, are provided along with failure analysis results. Mitigation techniques are also discussed. This analysis points towards an increased need for chip level fault-tolerance, especially error detecting and correcting codes, and hardware based checkpointing and retry.*

## 1. Introduction

Scaling of semiconductor devices has led to remarkable performance gains. At the same time, smaller transistor and interconnect features, lower supply voltage and increased clock frequency have contributed to higher error rates. Permanent, transient and intermittent faults are the main sources of errors in integrated circuits (IC).

Permanent faults occur due to irreversible physical changes. Shorts and opens are typical examples of such faults. Transients are most frequently generated by environmental conditions, like cosmic rays. Intermittent faults occur due to unstable or marginal hardware. Manufacturing residues may lead to such faults. Errors induced by transient and intermittent faults manifest similarly. However, two main criteria may be used to determine the source of an error. On one hand, errors induced by intermittent faults usually occur in bursts, at the same location, when the fault is activated. On the other, replacement of the offending part eliminates an intermittent fault, while transients cannot be fixed by repair. Additionally, intermittent faults may be activated or deactivated by temperature, voltage and frequency changes.

The effects of scaling on IC dependability have been extensively analyzed. For instance, both measurements and simulation were employed for determining soft error rate (SER) dependency on semiconductor technology [11, 18]. The impact of scaling on errors due to fluctuating minimum voltage was addressed in [1].
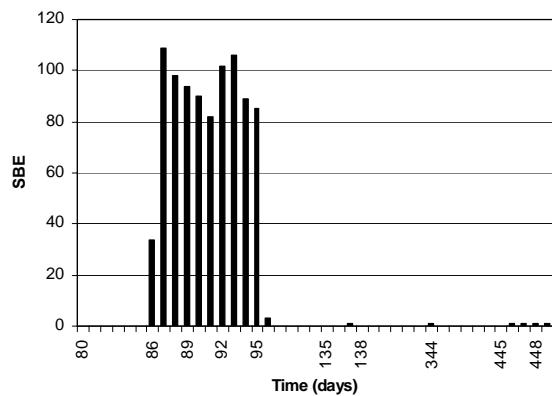
A variety of techniques have been devised for handling the errors induced by silicon faults. In particular, special attention has been paid to soft errors induced by high energy particles [2, 10]. However, specific behavior of the intermittent faults, especially error burstiness, was rarely considered.

This paper concentrates on the impact of intermittent faults on IC dependability, discusses the impact of scaling, and covers several competing mitigation techniques. An experiment for collecting field error data, from IT production servers, is briefly presented in Section 2. A typical signature, for memory errors induced by intermittent faults, is provided. Section 3 discusses the ultra-thin oxide breakdown and impact of scaling on this failure mode. Section 4 concentrates on timing violations which may manifest intermittently. Solutions for handling errors induced by silicon faults are discussed in Section 5. Section 6 concludes the paper.

## 2. Field error data

Determining the most frequent sources of errors, and their manifestation in real computing systems, is paramount for properly designing fault/error handling systems. To this end, 257 servers, produced by two manufacturers, were monitored. 310.7 server years worth of data was collected. Failure analysis was carried out, whenever feasible, for finding out the root cause of the errors.

Memory single-bit errors (SBE), induced by intermittent faults, were prevalent. 47.5% of the servers reported no errors, and 31.5% systems experienced one to five errors. However, a number of servers experienced large bursts of SBE: 5.8% servers logged 101 to 1000 SBE and 1.9% reported over 1000 errors.



**Fig. 1. Memory SBE signature, generated by an intermittent fault**

Analysis of the error logs showed that 6.2% of the memory subsystems were affected by intermittent faults. Fig. 1 shows a common memory SBE signature [9]. The x and y axes show the day of occurrence and the number of SBE, respectively. Failure analysis found that polymer residue led to an intermittent contact.

It is expected that microscopic manufacturing residues, tolerated by present technologies, will increase the rate of occurrence of the intermittent faults, as the device size shrinks down to a few nanometers.
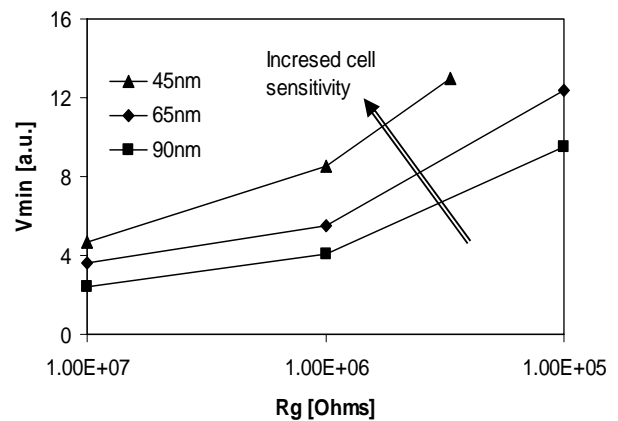
The monitoring system also logged SBE bursts experienced by the data signals of the processor bus, in the case of two servers. Neither service interruption nor silent data corruption (SDC) occurred, as the data path of the bus was protected by ECC. The processors connected to the bus reported SBE bursts, from 15 errors up to 7104 errors. Failure analysis revealed that solder joint intermittent contacts were the source of

the errors [9]. Additionally, physical and simulated fault injection experiments showed that similar faults on the processor control signals may lead to SDC.

The interconnect scaling is likely to increase the frequency of occurrence of these types of intermittent faults.

## 3. Oxide failures

Increased current leakage is becoming a serious concern as oxide layer thickness approaches 3 – 4 nm. Oxide breakdown starts with the tunnel injection of electrons, which induces microscopic defects. The higher leakage current generates thermal damage and the defect expands laterally, eventually leading to permanent breakdown.



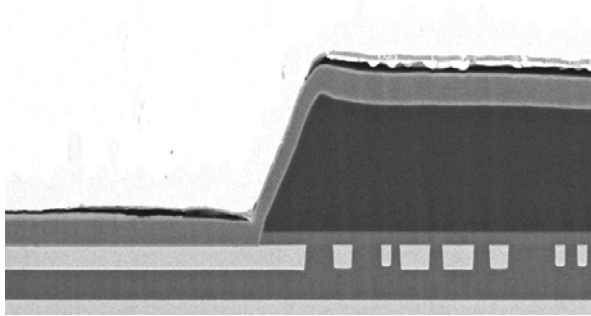**Fig. 2. Vmin as a function of gate oxide resistance and technology node [1]**

A new failure mechanism, specific to nanometer scale devices, is the soft breakdown (SBD). In this case the leakage current fluctuates, without inducing the thermal damage [19].

Effects of SBD have been already observed. For instance, erratic fluctuations of the minimum voltage (Vmin), in 90nm technology SRAM, were reported in [1]. In this case the measured Vmin varied from 0.55V to 0.8V, due to oxide leakage of a NMOS pull-down transistor. The intermittent behavior of the failure was attributed to gate oxide SBD. Better oxide and optimization of the pull-down transistor were needed to fix the problem.

Vmin sensitivity to gate current leakage is expected to increase with scaling. Fig. 2 shows Vmin dependency on gate oxide resistance, for three technology nodes [1].

## 4. Timing failures

Timing failures may occur due to propagation delays over the interconnect. For instance, the barrier layer material (BLM) delamination shown in Fig. 3 leads to higher resistance and, as a result, to timing violations [9]. Other failure mechanisms, like electromigration induced voids, increase the interconnect resistance, generating errors intermittently. Crosstalk delays may occur when adjacent signals switch in opposite directions [8]. Process, temperature, and voltage (PVT) variations tend to amplify this phenomenon. It is expected that crosstalk effects will increase with interconnect scaling and higher clock frequencies.



**Fig. 3. Increased resistance due to BLM delamination is a source of intermittent faults**

## 5. Mitigation techniques

Two main approaches are commonly employed for improving dependability of computing systems: fault avoidance and fault tolerance. Fault avoidance techniques have been extensively used for lowering the rate of occurrence of high-energy particle induced soft errors. Silicon on insulator (SOI) technology is an example [3, 4]. However, SOI does not have any significant impact on the rate of occurrence of the intermittent faults.

A new cost effective class of solutions, employing software fault-tolerance, has been proposed in the last few years. Multithreading based time redundancy techniques [15, 16] are able to handle rare events well, for instance particle induced upsets, but are less effective in the case of intermittent faults. Software only solutions experience significant performance penalties when large bursts of errors occur. Even

worse, the high frequency errors, specific to active intermittent faults, may lead to a near coincident fault scenario, i.e., a new error arrives before the handling of the previous one is completed, leading to the failure of the recovery process.

Data gathered in this study suggests that fault avoidance techniques should address not only particle induced errors, but intermittent faults as well. For example, the impact of crosstalk induced timing violations can be diminished by gate sizing [20].

Hardware implemented error handling techniques are likely to provide the best solutions for mitigating the effects of intermittent faults. The high speed of silicon logic makes hardware implementations well suited for detection and correction of errors occurring at a high rate. For instance, ECC deliver fast error detection and recovery [6, 7, 13]. Scrubbing techniques may be used in conjunction with ECC for avoiding accumulation of errors in memory arrays [17]. A novel hardware implemented technique, designed for improving register file reliability, is based on replication of narrow-width values [12]. This approach allows for error recovery as well, if both parity and replication are employed.

Hybrid solutions, which combine hardware error detection and recovery with software implemented failure prediction and resource reconfiguration, may improve dependability significantly. Such approaches are already employed by some high end servers. For instance, redundant instruction and execution units and hardware implemented machine state checkpointing, coupled with software controlled reconfiguration, is described in [5, 14].

## 6. Conclusions

Field data, laboratory measurements and simulation results suggest that intermittent faults represent a significant threat to dependability of nanocomputing devices. Manufacturing residuals, ultra-thin oxide degradation, and crosstalk induced delays are a few examples of such faults. Aggressive scaling and increased complexity are expected to lead to higher rates of occurrence of the intermittents, despite the extensive use of fault avoidance.

As a result, fault-tolerance techniques have to be widely employed. Software only solutions are too slow for effectively handling large bursts of errors. It is expected that complex integrated circuits, in general, and microprocessors, in particular, will provide extensive hardware fault-tolerance capabilities in the future. Software solutions will continue to significantly contribute to improving dependability of computing systems, especially by providing failure prediction and graceful performance degradation.

## Acknowledgement

## References

[1] M. Agostinelli et al, "Erratic fluctuations of SRAM cache Vmin at the 90nm process technology node", Proceedings of International Electron Devices Meeting, 2005, http:/www.intel.com/technology/silicon/micron.htm

[2] G. R. Agrawal, L. W. Massengill, K. Gulati, "A proposed SEU tolerant dynamic random access memory (DRAM) cell", IEEE Transactions on Nuclear Science, Vol. 41, No. 6, 1994, pp. 2035-2042.

[3] J. Baggio et al, "Neutron-induced SEU in bulk and SOI SRAMs in terrestrial environment", IEEE Reliability Physics Symposium, pp. 677-678, 2004.

[4] E. H. Cannon, D. R. Reinhardt, M. S. Gordon, P. S. Makowenskyj, "SRAM SER in 90, 130 and 180 nm bulk and SOI technologies", IEEE Reliability Physics Symposium, 2004, pp. 300-304.

[5] M. A. Check, T. HJ. Slegel, "Custom S/390 G5 and G6 microprocessors", IBM Journal of Research and Development, Vol. 43, No. 5/6, 1999, pp. 671-680.

[6] C. L. Chen, "Symbol error correcting codes for memory applications", Proceedings of FTCS, 1996, pp. 200- 207.

[7] C. L. Chen, "On double-byte error-correcting codes", IEEE Transactions on Information Theory, Vol. 45, No. 6, 1999, pp. 2207 – 2208.

[8] W. Y. Chen, S. K. Gupta, M. A. Breuer, "Test generation for crosstalk-induced faults: framework and computational results", Proceedings of Asian Test Symposium, 2000, pp. 305-310.

[9] C. Constantinescu, "Intermittent Faults in VLSI Circuits", IEEE Workshop on System Effects of Logic Soft Errors, 2006, www.selse.org

[10] K. Gulati, L. W. Massengill, G. R. Agrawal, "Single event mirroring and DRAM sense amplifier designs for improved single-event-upset performance", IEEE Transactions on Nuclear Science Vol. 41, No. 6, 1994, pp. 2026-2034.

[11] S. Hareland et al., "Impact of CMOS process scaling and SOI on the soft error rates of logic processes", IEEE Symposium on VLSI Technology, 2001, pp. 73-74.

[12] J. Hu, S. Wang, S. G. Ziavras, "In-Register Duplication: Exploiting Narrow-Width Value for Improving Register File Reliability", Proceedings of IEEE Dependable Systems and Networks Conference, 2006, pp. 281-290.

[13] Y. Katayama, S. Morioka, "One-shot Reed-Solomon decoding for high-performance dependable systems", Proceedings of IEEE Dependable Systems and Networks Conference, 2000, pp. 390 -399.

[14] P. J. Meaney, S. B. Swaney, P. N. Sanda, L. Spainhower, "IBM z990 soft error detection and recovery", IEEE Transactions on Device and Materials Reliability, Vol. 5, No. 3, 2005, pp. 419-427.

[15] F. Rashid, K. K. Saluja, P. Ramanathan, "Fault tolerance through re-execution in multiscalar architecture", Proceedings of IEEE Dependable Systems and Networks Conference, 2000, pp. 482 – 491.

[16] E. Rotenberg, "AR-SMT: A microarchitectural approach to fault tolerance in microprocessors", Proceedings of 29[th] FTCS Symposium, 1999, pp. 84-91.

[17] A. M. Saleh, J. J. Serrano, J. H. Patel, "Reliability of scrubbing recovery-techniques for memory systems", IEEE Transactions on Reliability, Vol. 39, No. 1,1990 , pp. 114 - 122.

[18] P. Shivakumar et al, "Modeling the effect of technology trends on the soft error rate of combinatorial logic", Proceedings of IEEE Dependable Systems and Networks Conference, 2002, pp. 389-398.

[19] J. H. Stathis, "Physical and predictive models of ultrathin oxide reliability in CMOS devices and circuits", IEEE Transactions on Device and Materials Reliability", Vol. 1, No. 1, 2001, pp. 43-99.

[20] T. Xiao, M. Marek-Sadowska, "Gate sizing to eliminate crosstalk induced timing violations", Proceedings of International Computer Design Conference, 2001, pp. 186-191.