

Lifetime Reliability Awareness for Microprocessors

Jayanth Srinivasan, Sarita V. Adve
University of Illinois at Urbana-Champaign
Department of Computer Science
{srinivsn,sadve}@cs.uiuc.edu,

Pradip Bose, Jude A. Rivers
IBM T.J. Watson Research Center
Yorktown Heights, NY
{pbose,jarivers}@us.ibm.com

Lifetime reliability due to wear-out related hard errors of processor components is emerging as a critical challenge in modern microprocessors. The steady processor performance increases seen over the last twenty years have been driven by aggressive scaling of CMOS devices. At the same time, scaling leads to reduced device feature sizes which results in lower processor lifetime reliability [4]. Device, manufacturing, and fabrication researchers have been aware of the lifetime reliability problem for many years and there exists a large body of research at the device level. On the other hand, there is a dearth of architectural lifetime reliability research as microarchitects have traditionally not viewed the subject as a problem. This poster highlights our work on microarchitectural lifetime reliability awareness.

As a first step towards understanding lifetime reliability from an architectural perspective, we proposed RAMP, a microarchitecture-level model that can be incorporated into a processor to dynamically track an estimate of lifetime reliability, accounting for the behavior of the executing application [3]. RAMP is based on state-of-the-art device level failure models obtained from researchers at IBM T.J. Watson Research Center. A brief overview of RAMP is presented in our poster.

We also integrated device scaling models in RAMP and quantified the impact of technology scaling on processor temperature and reliability in [4], showing that scaling has a significant and increasing effect on processor hard failure rates. For a POWER4-like processor running Spec2000 applications, they project an average increase of 316% in processor failure rates when scaling from 180nm to 65nm [4] (this implies that a processor that would take 10 years to fail at 180nm would fail in less than 3 years at 65nm). More importantly, they show that the rate of increase of failure rates also increases with scaling, clearly highlighting the challenges that will be imposed on lifetime reliability in the near future. These results are also presented in the poster.

Next, we explored microarchitectural techniques to enhance lifetime reliability and/or decrease design cost and time. Our work allows the system to view reliability as a resource and optimally distribute it across system components in an application aware fashion. **Dynamic reliability management (DRM)** is a technique where the processor adaptively responds to changing application behavior to maintain its lifetime reliability target, but possibly at a different performance [3]. DRM has several potential benefits – instead of assuming worst-case behavior for reliability qualification, DRM allows processors to be qualified for reliability at lower (but more likely) operating points, which are commensurately cheaper. Dynamic management serves as a backup to ensure that the target failure rate will

not be exceeded, although some performance may be lost. Conversely, if the processor is over-designed for reliability, DRM can be used to extract excess performance. Hence, DRM provides a hitherto unexplored design strategy for the system to consciously tradeoff performance and cost, to achieve the target reliability goals. We present an overview of DRM and our results in the poster.

We have also examined design time enhancements for processor reliability. In particular, we examine the reliability benefits of introducing structural redundancy during microarchitectural specification. Although redundancy has been commonly used for reliability, most previous work focused on redundancy at the processor granularity. Due to the large area overheads involved in duplicating entire processors, such redundancy does not provide a cost-effective reliability solution. In addition, most processor redundant systems require the replacement processor to be manually installed [1]. Structural redundancy addresses some of these shortcomings of processor redundancy by incurring less area overheads and allowing run-time processor reconfiguration for reliability.

Specifically, we examine two methods by which structural redundancy can be used for run-time reliability enhancement [2]. In the first case, **structural duplication (SD)**, certain redundant microarchitectural structures are added to the processor and designated as "spares". Spare structures can be turned on during the processor's lifetime when the original structure fails. We also propose a new run-time technique, **graceful processor degradation (GPD)**, which allows the processor to exploit existing microarchitectural redundancy for reliability. Modern processors have replicated structures that are used for increasing performance for some high parallelism applications. If a replicated structure fails in the course of a processor's lifetime, the processor can shut down the structure and still maintain functionality, thereby increasing lifetime. We compare structural duplication, graceful processor degradation, and a combination of the two techniques in our poster.

References

- [1] L. Spainhower and T. A. Gregg. Ibm s/390 parallel enterprise server g5 fault tolerance: A historical perspective. In *IBM Journal of Research and Development*, September/November 1999.
- [2] J. Srinivasan et al. Early Stage Microarchitectural Design for Lifetime Reliability. Submitted for publication.
- [3] J. Srinivasan et al. The Case for Lifetime Reliability-Aware Microprocessors. In *ISCA31*, June 2004.
- [4] J. Srinivasan et al. The Impact of Technology Scaling on Lifetime Reliability. In *DSN04*, June 2004.