

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/318728082>

Quality-Configurable Memory Hierarchy Through Approximation

Conference Paper · October 2017

CITATIONS

0

READS

95

3 authors, including:



Amir M. Rahmani

UC Irvine (USA) & TU Wien (Austria)

152 PUBLICATIONS **597** CITATIONS

[SEE PROFILE](#)



Nikil Dutt

University of California, Irvine

623 PUBLICATIONS **11,459** CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Soft Errors [View project](#)



IoCT-CARE: Internet of Cognitive Things for Personalized Healthcare [View project](#)

All content following this page was uploaded by [Amir M. Rahmani](#) on 28 July 2017.

The user has requested enhancement of the downloaded file.

and/or low write endurance. To reduce energy consumption of these memories and improve the lifetime of NVMs, approximate storage techniques sacrifice data integrity by reducing supply voltage in SRAM [11] [15] and refresh rate in (e)DRAM [16] [3] [9] and by relaxing or skipping read/write operation in NVMs [25] [17] [4] [2] [1] [22] [7] [14] [6].

These objectives are achieved by a number of general strategies, namely: (1) precision scaling [20] [18], (2) approximating load operations [5] [19] [21], (3) skipping store operations [5] [17] [13], (4) using faulty or unprotected memory substrate [22] [4] [10] [15] [7] [14] [6] [25], (5) tweaking technology-dependent reliability-energy knobs [1] [11] [16] [3] [2] [4] [9].

3 EXEMPLARS OF QUALITY-CONFIGURABLE MEMORIES

3.1 Relaxed Cache

Relaxed cache [11] is a technique for saving leakage energy in SRAM-based caches. In Relaxed Cache a subset of the cache ways are protected (using error detection/correction and/or proper supply voltage settings). Other cache ways (called relaxed ways) operate with a lower supply voltage and therefore may become faulty. The programmer identifies noncritical data objects using suitable annotations, and the virtual addresses of these objects are kept in a table. On a cache miss, if the missed data block is noncritical, a victim block is selected from the relaxed ways. A critical data item is always stored in a fault-free block.

To control the degree of approximation, any cache block that has more than a certain acceptable number of faulty bits (called AFB) is disabled and power-gated. AFB along with VDD are two knobs used in this work to determine the effective capacity of the cache and also bound the overall error exposed to software. The experimental results show that Relaxed Cache save significant amount of cache leakage energy (up to 74%) while still generating acceptable quality results.

3.2 QuARK

QuARK [1] is a hardware/software approach for trading reliability of STT-MRAM caches for energy savings in the on-chip memory hierarchy of systems running approximate applications. QuARK utilizes fine-grained actuation knobs to efficiently control reliability-energy trade-offs for individual accesses of concurrently running applications.

Compared to the related work on SRAM caches [11] that uses cache-way-level knobs, QuARK presents a more fine-grained actuation capability enabling it to offer the following advantages: i) the knob actuations do not affect any other cache block unlike the actuation in [11] which requires flushing all the affected blocks, ii) multiple applications with different degrees of reliability can share the same cache without affecting each other's guaranteed level of reliability, and iii) the reliability level requested for a piece of data can be changed at runtime, if the quality of output is not satisfactory. The simulation results demonstrated up to 40% energy savings over a fully-protected STT-MRAM L2 cache, with acceptable quality loss.

4 CONCLUSIONS AND FUTURE DIRECTIONS

In this paper, we presented state-of-the-art approximate memory management approaches and discussed opportunities for energy

savings in the on-chip memory hierarchy. We believe the gains from approximation can be maximized with cross-layered techniques that can match the possibilities in one layer with opportunities in the other layers. However, to realize such opportunities, we need to cope with open challenges such as: i) how to develop programming models to capture the resiliency of data objects of a program, ii) how to automate classification of data objects based on their resiliency to memory errors, iii) how to achieve a holistic approach to guide approximation across the entire memory hierarchy, iv) how to orchestrate memory approximation knobs at different levels of abstraction, and v) how to control the approximation process of the memory subsystem to bound errors propagated to applications.

ACKNOWLEDGMENT

We acknowledge financial support by the Marie Curie Actions of the European Union's H2020 Programme.

REFERENCES

- [1] A. M. Monazzah et al. 2017. QuARK: Quality-configurable Approximate STT-MRAM Cache by Fine-grained Tuning of Reliability-Energy Knobs. In *Proc. of ISLPED*.
- [2] A. Ranjan et al. 2015. Approximate Storage for Energy Efficient Spintronic Memories. In *Proc. of DAC*.
- [3] A. Raha et al. 2017. Quality Configurable Approximate DRAM. *IEEE Trans. Comput.* (2017).
- [4] A. Sampson et al. 2013. Approximate Storage in Solid-state Memories. In *Proc. of MICRO*.
- [5] B. Thwaites et al. 2014. Rollback-free Value Prediction with Approximate Loads. In *Proc. of PACT*.
- [6] D. Jevdjic et al. 2017. Approximate Storage of Compressed and Encrypted Videos. In *Proc. of ASPLOS*.
- [7] F. Sampaio et al. 2015. Approximation-aware Multi-Level Cells STT-RAM Cache Architecture. In *Proc. of CASES*.
- [8] G. P. Arumugam et al. 2015. Novel Inexact Memory Aware Algorithm Co-design for Energy Efficient Computation: Algorithmic Principles. In *DATE*.
- [9] K. Cho et al. 2014. eDRAM-based Tiered-Reliability Memory with Applications to Low-power Frame Buffers. In *Proc. of ISLPED*.
- [10] K. Lee et al. 2006. Mitigating Soft Error Failures for Multimedia Applications by Selective Data Protection. In *Proc. of CASES*.
- [11] M. Shoushtari et al. 2015. Exploiting Partially-Forgetful Memories for Approximate Computing. *IEEE Embedded Systems Letters* (2015).
- [12] N. Dutt et al. 2014. Multi-Layer Memory Resiliency. In *Proc. of DAC*.
- [13] O. Kislal et al. 2016. Cache-Aware Approximate Computing for Decision Tree Learning. In *Proc. of IPDPSW*.
- [14] Q. Guo et al. 2016. High-Density Image Storage Using Approximate Memory Cells. In *Proc. of ASPLOS*.
- [15] S. Ganapathy et al. 2015. Mitigating the Impact of Faults in Unreliable Memories for Error-resilient Applications. In *Proc. of DAC*.
- [16] S. Liu et al. 2011. Flicker: Saving DRAM Refresh-power Through Critical Data Partitioning. In *Proc. of ASPLOS*.
- [17] Y. Fang et al. 2012. SoftPCM: Enhancing Energy Efficiency and Lifetime of Phase Change Memory in Video Applications via Approximate Write. In *Proc. of ATS*.
- [18] Y. Tian et al. 2015. ApproxMA: Approximate Memory Access for Dynamic Precision Scaling. In *Proc. of GLSVLSI*.
- [19] J. S. Miguel et al. 2014. Load Value Approximation. In *IEEE/ACM International Symposium on Microarchitecture*.
- [20] J. S. Miguel et al. 2015. DoppelgÄnger: A Cache for Approximate Computing. In *Proc. of MICRO*.
- [21] J. S. Miguel et al. 2016. The Bunker Cache for Spatio-value Approximation. In *Proc. of MICRO*.
- [22] F. Oboril, A. Shirvanian, and M. Tahoori. 2016. Fault Tolerant Approximate Computing using Emerging Non-volatile Spintronic Memories. In *Proc. of VTS*.
- [23] Martin Rinard. 2013. Parallel Synchronization-Free Approximate Data Structure Construction. In *5th USENIX Workshop on Hot Topics in Parallelism*.
- [24] Majid Shoushtari and Nikil Dutt. 2017. *A Survey of Techniques for Approximate Memory Management*. Technical Report CECS-TR-17-03. Center for Embedded and Cyber-physical Systems, University of California, Irvine.
- [25] X. Xu and H. H. Huang. 2015. Exploring Data-Level Error Tolerance in High-Performance Solid-State Drives. *IEEE Transactions on Reliability* (2015).