

Nonblocking DRAM Refresh

Kate Nguyen, Kehan Lyu, Xianze Meng

Department of Computer Science

Virginia Tech

Blacksburg, Virginia

katevy@vt.edu, kehan@vt.edu, xianze@vt.edu

Vilas Sridharan

RAS Architecture

Advanced Micro Devices, Inc

Boxborough, Massachusetts

vilas.sridharan@amd.com

Xun Jian

Department of Computer Science

Virginia Tech

Blacksburg, Virginia

xunj@vt.edu

Abstract—Since its invention half a century ago, DRAM (Dynamic RAM) has required dynamic refresh operations that block read accesses to refreshing data; this fundamental behavior gave DRAM its name. In contrast, DRAM's close relative - SRAM (static RAM) - can statically re-enforce charge in the background without blocking read accesses at the cost of more expensive circuit structure. Nonblocking DRAM Refresh blurs this fundamental distinction between DRAM and SRAM at the system level to enable the best of both worlds - allowing read accesses to refreshing data in DRAM while preserving DRAM's low-cost circuit structure.

DRAM is arguably the most dominant type of memory technology today. Due to emerging trends such as 3D die-stacked DRAM in processor packages (e.g., HBM) and embedded DRAM on processor dies (e.g., to implement large LLCs), DRAM plays critical role in computing. However, since its inception a half a century ago, DRAM has had an inherent physical characteristic that increases latency compared to its close relative - Static RAM (SRAM). While DRAM and SRAM are both volatile, DRAM requires dynamic/active refresh operations that stall accesses to refreshing data; in comparison, SRAM statically re-enforces electric charge in the background without stalling accesses but requires a more expensive circuit structure (e.g., 6T1C versus 1T1C for DRAM).

Refresh-induced stalls in DRAM reduce system performance. This has been a universal feature for systems using any DRAM memory. The performance overhead also grows as DRAM density increases; denser DRAM devices have more cells to refresh and, therefore, require either longer or more frequent refresh operations. Figure 1 shows DRAM refresh latency versus density for DDRx DRAM.

Nonblocking DRAM refresh eliminates this fundamental difference between DRAM and SRAM at the system level by refreshing DRAM in the background without stalling read requests to refreshing memory blocks. In a conventional DRAM refresh scheme, a single refresh operation refreshes all of the data in a memory block simultaneously. Under Nonblocking DRAM Refresh, on the other hand, a single refresh operation refreshes only a portion of the data in a memory block. A system implementing nonblocking refresh equips each memory block with enough redundant data (e.g., parity) to reconstruct the blocks refreshing data when it

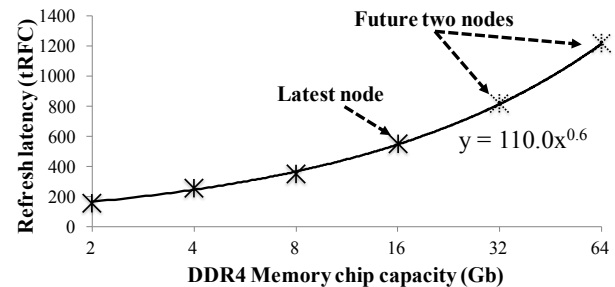


Figure 1. Historical (solid markers) and projected (dashed markers) refresh latency in nanoseconds for DDR4 DRAM chips.

is accessed, allowing the memory controller to issue reads to refreshing blocks. The memory controller issues refresh operations more frequently to keep the overall refresh time in memory constant.

The primary cost of implementing nonblocking refresh is equipping each memory block with redundant data to compute the refreshing/unreadable data in refreshing memory blocks. However, many existing memory systems already contain substantial redundant data to protect against hardware failures. We observe that these redundant data are budgeted to protect against worst-case hardware failure scenarios. Therefore, the redundant data are under-utilized in the common case when hardware faults are absent. A memory block that has under-utilized redundant data has *memory error resilience slack*. Memory architectures containing redundant data for hardware failure protection can leverage their memory error resilience slack to implement nonblocking refresh with little additional overhead.

Nonblocking refresh can be applied across all DRAM-based memories (e.g., off-package DRAM, in-package DRAM, and embedded DRAM); however, the microarchitectural details (e.g., how to refresh only a portion of data in a memory block at a time) vary depending on the particular DRAM technology. This article explores the microarchitectural details for making DRAM refresh nonblocking for server memory systems built from DDRx DRAM. We end the article by also briefly describing how to enable Nonblocking DRAM Refresh for other classes of DRAM memories.

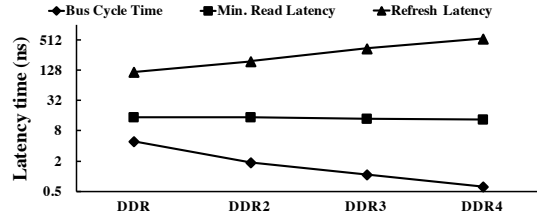


Figure 2. Historical scaling of memory latencies [1], [2].

IMPACT OF DRAM REFRESH ON SYSTEM PERFORMANCE

Each DRAM cell stores one bit of data as electric charge. However, the charge leaks over time; a DRAM cell loses its data if it loses its charge. To refresh many DRAM cells in memory sufficiently quickly, many DRAM cells are refreshed together; this causes large regions in memory to be unreadable during each refresh operation. In DDRx DRAM, the entire DRAM chip is unreadable when the chip performs refresh. As the DRAM density increases with newer generations of DRAM, refresh latency also grows because each refresh operation has to refresh more cells. In contrast, other memory latencies have remained steady or decreased across generations. Figure 2 shows historical DRAM latency scaling; it shows that other DRAM latencies, such as bus cycle time and minimum read latency, have improved while refresh latency has increased. If these trends continue, refresh latency may become a determining factor in overall memory system performance.

The inability to read data from refreshing chips stalls program execution. To mitigate the resultant performance overhead, many recent works have proposed skipping half or more of the refresh operations to improve performance [3], [4], [5], [6]. For example, RAIDR [4] profiles the retention time of DRAM cells and skips refresh operations to memory cells with long retention time. However, skipping refresh reduces the average amount of charge stored in DRAM cells and, therefore, can increase DRAM vulnerability to read disturb errors [7].

MAKING DRAM REFRESH NONBLOCKING FOR SERVER MEMORY SYSTEMS

A. Background on Server Memory Systems

Modern server memory architectures operate groups of DRAM chips in lockstep. A group of DRAM chips that operate in lockstep is called a *rank*. For example, when the processor issues a read request for a 64B memory block to a rank, every chip in the rank replies with a part (e.g., 4B or 8B) of the memory block in lockstep. Similarly, when the processor issues a refresh command to a rank, all chips in the rank refresh in lockstep.

In server memory systems, each memory block may contain redundant data for hardware failure protection. These redundant data can correct up to complete memory chips

failure(s) in a rank. The ratio of redundant data to program data in each block ranges from 12.5% to 40.6%.

B. Nonblocking DRAM Refresh Operations

Each Nonblocking DRAM Refresh operation refreshes a limited amount of data in a block at a time; this makes it possible for the server processor to use each memory block's redundant data to compute the block's small amount of data that are currently unreadable due to refresh (see detail later). In comparison, each conventional refresh operation refreshes entire blocks at the same time because it refreshes all chips in a rank in lockstep.

To limit the amount of data that is refreshing in a memory block at a time, the memory controller issues a refresh command to only a few DRAM chips in a rank at a time. Limiting the amount of data that is refreshing in a memory block at a time enables the memory controller to use the memory blocks redundant data to reconstruct the small amount of currently inaccessible data in the block. Because Nonblocking Refresh operations do not stall read requests, the memory controller can compensate for refreshing only some of the data in a block at a time by issuing refresh commands more frequently. In comparison, conventional refresh operations, which block read requests, can only refresh each rank infrequently to avoid excessively stalling read requests. A server system, for instance, can issue a refresh operation to the next subgroup of chips in a rank back-to-back in round robin fashion as soon as the previous subgroup finishes refresh

C. Reading from Refreshing Memory Blocks

Modern server systems equip each memory block with sufficient redundant data to correct a complete memory chip failure. The redundant data are under-utilized in the common case when hardware faults are absent. Figure 3 quantifies the expected fraction of non-faulty memory pages over time based on a large-scale field study [8], assuming eight ranks per channel and 18 chips per rank. On average across seven years of operations, 97% of memory pages are not affected by any fault. While only a small fraction of memory experience faults, server systems protect all memory blocks with uniform redundant data because memory faults are stochastic events, whose time and location are difficult to predict.

A memory block that has under-utilized redundant data has *memory error resilience slack*. Nonblocking DRAM Refresh uses memory error resilience slack to compute the unreadable data residing in refreshing chips without compromising reliability. When the processor fetches a block from a rank currently performing Nonblocking DRAM Refresh, the part of its data that reside in refreshing chips will be missing in the fetched block. The processor can use the fetched block's redundant data to compute the missing data. Figure 4 illustrates how to read from a rank that is



Figure 3. Expected fraction of memory pages that have not yet been affected by any hardware fault as a function of time.

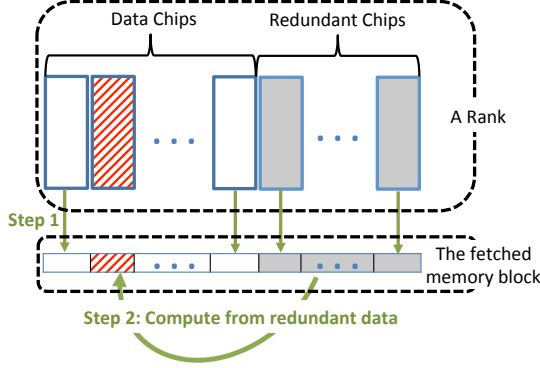


Figure 4. Reading a block from a rank performing Nonblocking DRAM Refresh. Red stripes represent data currently unreadable due to refresh.

currently performing Nonblocking Refresh. Computing the missing data is fast because the processor knows which memory chip(s) are refreshing and, therefore, which part of the block’s data is missing. This is unlike regular error correction, where the processor needs to first determine where is the missing data before it can compute the correct value at that location; the majority of the latency incurred by error correction is due to locating the error [9].

A system implementing Nonblocking Refresh must maintain the same level of error detection and correction as a system with conventional refresh. In a server memory system, each memory block contains some redundant data for error detection and some redundant data for error correction. To preserve error detection strength, a system with Nonblocking Refresh uses the same amount of redundant data to detect hardware errors for each read request as baseline systems. Nonblocking Refresh only uses the remaining redundant data in the memory block to compute missing data due to refresh. To preserve error correction strength, the memory controller computes data missing due to refresh only when the memory controller does not detect any hardware errors. When the memory controller does detect hardware error(s) for a read request, the processor waits for the rank to finish its in-flight refresh and then fetches the previously-inaccessible data directly from DRAM. As all the data from the block are now available, the memory controller performs error correction exactly the same way as in a conventional system

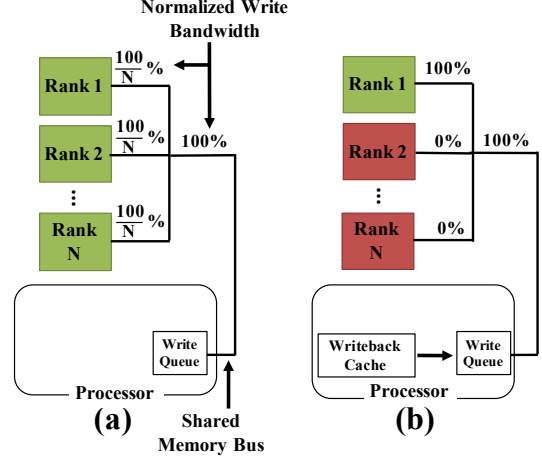


Figure 5. Write distribution in (a) conventional and (b) proposed memory systems. Green ranks are not refreshing and, therefore, writable; red ranks are refreshing and, therefore, not writable.

and, therefore, preserves the same error correction strength.

D. Write Requests During Nonblocking Refresh

Write requests to refreshing blocks still need to wait for refresh to complete; this is challenging for Nonblocking Refresh as it must refresh each memory block more frequently. Because write requests are not on the critical path of program execution, write latency has low impact on performance; as such, the processor may simply buffer write requests until the corresponding rank is not refreshing by adding a small (e.g., 32KB) write-buffer cache to the memory controller.

Because all ranks in the same channel share the same memory bus, the processor can only write to one rank at a time in a channel. Therefore, total write bandwidth in a channel is divided across all the ranks in the channel as illustrated in Figure 5A. Second, all ranks in a channel receive fairly even number of writes because logically adjacent memory pages are often interleaved across ranks to maximize rank-level parallelism. As such, the processor can re-order write requests to concentrate each channel’s write bandwidth to one rank at a time as shown in Figure 5B. This maintains the same channel-level write bandwidth while allowing the remaining ranks in the channel to perform Nonblocking Refresh.

MAKING DRAM REFRESH NONBLOCKING FOR OTHER TYPES OF MEMORY SYSTEMS

Nonblocking DRAM Refresh can apply to all DRAM-based memory systems. For example, desktop/laptop memory systems use the same rank architecture as server memory systems, except that they do not contain redundant data chips. As such, the same implementation of Nonblocking DRAM Refresh described for server memory systems is applicable to desktop/laptop memory systems after adding a redundant chip per rank to these memory systems.

Nonblocking DRAM Refresh is also applicable to memory systems that access only one DRAM chip per memory request, such as High Bandwidth Memory (HBM), graphics DDR (e.g., GDDRx), and smartphone memory (e.g., LPDDRx), because the internal organization within each DRAM die mirrors a memory channel's organization. For example, an HBM die contains multiple banks that share a common data bus [10], just like DDRx contains multiple ranks in a channel sharing a common data bus. There are also multiple sub-arrays per bank just like there multiple chips per rank [10]. In addition, each memory block is spread across multiple sub-arrays in one bank of a DRAM die, just like how a memory block is spread across a rank [10]. As such, HBM devices can implement Nonblocking DRAM Refresh by refreshing a portion of the sub-arrays in a bank at a time and adding redundant sub-arrays to each bank to compute the unreadable data in refreshing sub-arrays.

RESULTS

Using the Gem5 architectural simulator and Ramulator DRAM simulator, we simulated various NASA Parallel, Parsec, and SPEC2006 benchmarks. We simulated 16-core processor with four memory channels and four ranks per channel.

Figure 6 shows the average performance benefit when applying Nonblocking DRAM Refresh to Intel/AMD and IBM Power8 server memory systems. For 16Gb DRAM chips, DRAM Refresh provides 13% and 17% average performance improvement for Intel/AMD and IBM server systems, respectively. Nonblocking DRAM Refresh may become even more beneficial for future DRAM chips with longer refresh latency as DRAM density increases; for our extrapolated 32Gb DRAM refresh latency, the average performance improvement increases to 21% and 35% for Intel/AMD and IBM server systems, respectively. The performance benefits for Intel/AMD server systems are lower because their amount of redundant data is lower; as such, to allow a refreshing memory block to be still calculable from its redundant data, each Nonblocking DRAM Refresh operation can only refresh one chip per rank at time, which provides low DRAM refresh throughput. As such, some conventional full-rank refresh operations that block read requests are still occasionally needed to achieve the same DRAM refresh throughput as conventional systems, resulting in reduced performance benefits.

CONCLUSIONS

Since its invention half a century ago, DRAM has required dynamic refresh operations that block read accesses to refreshing data; this fundamental behavior gave DRAM its name. Nonblocking DRAM Refresh improves upon this fundamental DRAM behavior by allowing read accesses to refreshing memory blocks. It operates at the micro-architecture level; it refreshes only a limited amount of data

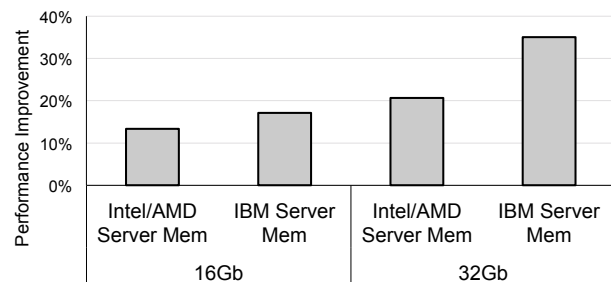


Figure 6. Average performance benefits of Nonblocking DRAM Refresh.

in a memory block at a time and uses redundant data in the block to compute the block's refreshing/unreadable data to complete read requests. The primary cost of Nonblocking DRAM Refresh is requiring some redundant data in each memory block. In systems containing redundant data for hardware failure protection, Nonblocking DRAM Refresh can eliminate this primary cost by safely leveraging memory error resilience slack to calculate the unreadable data in refreshing memory blocks. In this article, we demonstrate in detail such a memory-overhead-free Nonblocking DRAM Refresh design for server memory systems.

DRAM density scaling has been slowing down. An important reason is that increasing DRAM density often requires reducing DRAM cell size, which in turn reduces DRAM retention time; reduced DRAM retention time requires more frequent DRAM refresh operations, which incur higher performance overhead. Also exacerbating the performance overhead is the accompanying fact that increasing DRAM density means more cells need to be refreshed. By effectively tackling DRAM refresh overheads, Nonblocking DRAM Refresh can help sustain density scaling for all DRAM-based memories, as the high-level idea of Nonblocking DRAM Refresh is applicable to all DRAM-based memories that are accessed at the memory block granularity. These include DDRx DRAMs used in servers and personal computers, GDDRx DRAMs used in GPUs, mobile DRAMs used in smartphones, in-package 3D die-stacked DRAMs used in HPC systems, and embedded DRAMs used in LLCs within processors.

REFERENCES

- [1] M. T. Inc., "Speed vs. latency: why cas latency isn't an accurate measure of memory performance," 2015. <https://pics.crucial.com/wcsstore/CrucialSAS/pdf/en-us-c3-whitepaper-speed-vs-latency-letter.pdf>.
- [2] Micron, "8gb: x4, x8, x16 ddr4 sdram." https://www.micron.com/media/documents/products/data-sheet/dram/ddr4/8gb_ddr4_sdram.pdf.
- [3] A. Sampson, W. Dietl, E. Fortuna, D. Gnanaprasagam, L. Ceze, and D. Grossman, "Enerj: Approximate data types for safe and general low-power computation," *SIGPLAN Not.*, vol. 46, pp. 164–174, June 2011.

- [4] J. Liu, B. Jaiyen, R. Veras, and O. Mutlu, "Raidr: Retention-aware intelligent dram refresh," in *Computer Architecture (ISCA), 2012 39th Annual International Symposium on*, pp. 1–12, June 2012.
- [5] I. Bhati, Z. Chishti, S. L. Lu, and B. Jacob, "Flexible auto-refresh: Enabling scalable and energy-efficient dram refresh reductions," in *2015 ACM/IEEE 42nd Annual International Symposium on Computer Architecture (ISCA)*, pp. 235–246, June 2015.
- [6] M. Patel, J. S. Kim, and O. Mutlu, "The reach profiler (reaper): Enabling the mitigation of dram retention failures via profiling at aggressive conditions," in *Proceedings of the 44th Annual International Symposium on Computer Architecture, ISCA '17*, (New York, NY, USA), pp. 255–268, ACM, 2017.
- [7] Y. Kim, R. Daly, J. Kim, C. Fallin, J. H. Lee, D. Lee, C. Wilkerson, K. Lai, and O. Mutlu, "Flipping bits in memory without accessing them: An experimental study of dram disturbance errors," in *2014 ACM/IEEE 41st International Symposium on Computer Architecture (ISCA)*, pp. 361–372, June 2014.
- [8] V. Sridharan, J. Stearley, N. DeBardleben, S. Blanchard, and S. Gurumurthi, "Feng shui of supercomputer memory: Positional effects in dram and sram faults," in *Proceedings of SC13: International Conference for High Performance Computing, Networking, Storage and Analysis, SC '13*, (New York, NY, USA), pp. 22:1–22:11, ACM, 2013.
- [9] A. Kumar and S. Sawitzki, "High-throughput and low-power architectures for reed solomon decoder," in *Conference Record of the Thirty-Ninth Asilomar Conference on Signals, Systems and Computers, 2005.*, pp. 990–994, October 2005.
- [10] B. Giridhar, M. Cieslak, D. Duggal, R. Dreslinski, H. M. Chen, R. Patti, B. Hold, C. Chakrabarti, T. Mudge, and D. Blaauw, "Exploring dram organizations for energy-efficient and resilient exascale memories," in *2013 SC - International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*, pp. 1–12, Nov 2013.