

ETH 263-2210-00L COMPUTER ARCHITECTURE, FALL 2019

HW 1: DRAM REFRESH

Instructor: Prof. Onur Mutlu

TAs: Mohammed Alser, Rahul Bera, Geraldo Francisco De Oliveira Junior, Can Firtina, Juan Gomez Luna, Jawad Haj-Yahya, Hasan Hassan, Konstantinos Kanellopoulos, Jeremie Kim, Nika Mansouri Ghiasi, Lois Orosa Nogueira, Jisung Park, Minesh Hamenbhai Patel, Abdullah Giray Yaglikci

Given: Wednesday, Sep 25, 2019

Due: **Thursday, Oct 10, 2019**

- **Handin - Critical Paper Reviews (1).** You need to submit your reviews to <https://safari.ethz.ch/review/architecture19/>. Please, check your inbox, you should have received an email with the password you should use to login. If you didn't receive any email, contact comparch@lists.ethz.ch. In the first page after login, you should click in "Architecture - Fall 2019 Home", and then go to "any submitted paper" to see the list of papers.
- **Handin - Questions (2-5).** You should upload your answers to the Moodle Platform (<https://moodle-app2.1et.ethz.ch/mod/assign/view.php?id=382814>) as a single PDF file.

1. Critical Paper Reviews [300 points]

Please read the guidelines for reviewing papers and check the sample reviews. You may access them by *simply clicking on the QR codes below or scanning them*. We will give out extra credit that is worth 0.5% of your total grade for each good review.



Guidelines



Sample reviews

Write an approximately one-page critical review for each of the following papers. A review with bullet point style is more appreciated. Try not to use very long sentences and paragraphs. Keep your writing and sentences simple. Make your points bullet by bullet, as much as possible.

- Moscibroda and Mutlu, "Memory Performance Attacks: Denial of Memory Service in Multi-Core Systems," in Proceedings of the USENIX Security, 2007. https://people.inf.ethz.ch/omutlu/pub/mph_usenix_security07.pdf
- Liu et al., "RAIDR: Retention-Aware Intelligent DRAM Refresh," in Proceedings of the International Symposium on Computer Architecture, 2012. https://people.inf.ethz.ch/omutlu/pub/raidr-dram-refresh_isca12.pdf
- Kim et al., "Flipping Bits in Memory Without Accessing Them: An Experimental Study of DRAM Disturbance Errors", in Proceedings of the International Symposium on Computer Architecture, 2014. https://people.inf.ethz.ch/omutlu/pub/dram-row-hammer_isca14.pdf

2. Main Memory [150 points]

Answer the following questions for a machine that has a 4 GB DRAM main memory system and each row is refreshed every 64 ms.

- (a) During standalone runs of two applications (A and B) on the machine, you observe that application A spends a significantly larger fraction of cycles stalling for memory than application B does while both applications have a similar number of memory requests. What might be the reasons for this?

- (b) Application A also consumes a much larger amount of memory energy than application B does. What might be the reasons for this?

- (c) When applications A and B run together on the machine, application A's performance significantly degrades, while application B's performance does not degrade as much. Why might this happen?

- (d) The designer decides to use a smarter policy to refresh the memory. A row is refreshed only if it has not been accessed in the past 64 ms. Do you think this is a good idea? Why or why not?

- (e) When this new refresh policy is applied, the refresh energy consumption drops significantly during a run of application B. In contrast, during a run of application A, the refresh energy consumption reduces only slightly. Why might this happen?

3. DRAM Refresh - Utilization [150 points]

A memory system has four channels, and each channel has two ranks of DRAM chips. Each memory channel is controlled by a separate memory controller. Each rank of DRAM contains eight banks. A bank contains 32K rows. Each row in one bank is 8KB. The minimum retention time among all DRAM rows in the system is 64 ms. In order to ensure that no data is lost, every DRAM row is refreshed once per 64 ms. Every DRAM row refresh is initiated by a command from the memory controller which occupies the command bus on the associated memory channel for 5 ns and the associated bank for 40 ns. Let us consider a 1.024 second span of time.

We define *utilization* (of a resource such as a bus or a memory bank) as the fraction of total time for which a resource is occupied by a refresh command.

For each calculation in this section, you may leave your answer in *simplified* form in terms of powers of 2 and powers of 10.

- (a) How many refreshes are performed by the memory controllers during the 1.024 second period in total across all four memory channels?

- (b) What command bus utilization, across all memory channels, is directly caused by DRAM refreshes?

- (c) What data bus utilization, across all memory channels, is directly caused by DRAM refreshes?

- (d) What bank utilization (on average across all banks) is directly caused by DRAM refreshes?

- (e) The system designer wishes to reduce the overhead of DRAM refreshes in order to improve system performance and reduce the energy spent in DRAM. A key observation is that not all rows in the DRAM chips need to be refreshed every 64 ms. In fact, rows need to be refreshed only at the following intervals in this particular system:

Required Refresh Rate	Number of Rows
64 ms	2^5
128 ms	2^9
256 ms	all other rows

Given this distribution, if all rows are refreshed only as frequently as required to maintain their data, how many refreshes are performed by the memory controllers during the 1.024 second period in total across all four memory channels?

What command bus utilization (as a fraction of total time) is caused by DRAM refreshes in this case?

(f) What DRAM data bus utilization is caused by DRAM refreshes in this case?

(g) What bank utilization (on average across all banks) is caused by DRAM refreshes in this case?

(h) The system designer wants to achieve this reduction in refresh overhead by refreshing rows less frequently when they need less frequent refreshes. In order to implement this improvement, the system needs to track every row's required refresh rate. What is the minimum number of bits of storage required to track this information?

(i) Assume that the system designer implements an approximate mechanism to reduce refresh rate using Bloom filters, as we discussed in class. One Bloom filter is used to represent the set of all rows which require a 64 ms refresh rate, and another Bloom filter is used to track rows which require a 128 ms refresh rate. The system designer modifies the memory controller's refresh logic so that on every potential refresh of a row (every 64 ms), it probes both Bloom filters. If either of the Bloom filter probes results in a "hit" for the row address, and if the row has not been refreshed in the most recent length of time for the refresh rate associated with that Bloom filter, then the row is refreshed. (If a row address hits in both Bloom filters, the more frequent refresh rate wins.) Any row that does not hit in either Bloom filter is refreshed at the default rate of once per 256 ms.

The false-positive rates for the two Bloom filters are as follows:

Refresh Rate Bin	False Positive Rate
64 ms	2^{-20}
128 ms	2^{-8}

The distribution of required row refresh rates specified in part (e) still applies.

How many refreshes are performed by the memory controllers during the 1.024 second period in total across all four memory channels?

What command bus utilization results from this refresh scheme?

What data bus utilization results from this refresh scheme?

What bank utilization (on average across all banks) results from this refresh scheme?

4. DRAM Refresh - Energy [150 points]

A new supercomputer has a DRAM-based memory system with the following configuration:

- The total capacity is 1 ExaByte (EB).
- The DRAM row size is 8 KiloByte (KB).
- The minimum retention time among all DRAM rows in the system is 64 ms. In order to ensure that no data is lost, every DRAM row is refreshed once every 64 ms. (Note: For each calculation in this question, you may leave your answer in simplified form in terms of powers of 2 and powers of 10.)

(a) How many DRAM rows does the memory system have?

(b) According to the specs of this hypothetical DRAM device, the memory controller should issue a refresh command (REF) at every 7.8 μ s. How many rows is refreshed when a REF command is issued?

(c) What is the power consumption of DRAM refresh? (Hint: you will need to figure out how much current the DRAM device draws during refresh operations. You can find useful information in the technical note by Micron¹. Use the current (IDD) numbers specified in Micron's datasheet².

Clearly state all the assumptions and show how you derive the power numbers. You are welcome to use other datasheets as well. Make sure you specify how you obtain the power numbers and show your calculations and thought process.)

(d) What is the total energy consumption of DRAM refresh during a refresh cycle? And during a day?

¹Micron, "TN-47-04: Calculating Memory System Power for DDR2 Introduction"

² Micron, "1Gb: x4, x8, x16 DDR3 SDRAM Features"

5. VRL: Variable Refresh Latency [150 points]

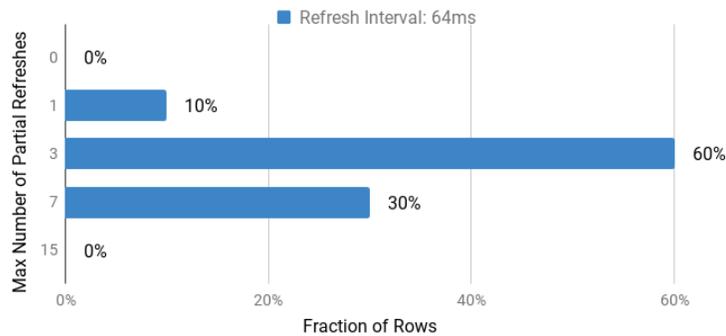
In this question, you are asked to evaluate Variable Refresh Latency³, which is based on two key observations:

- First, a cell's charge reaches 95% of the maximum charge level in 60% of the nominal latency value during a refresh operation. In other words, the last 40% of the refresh latency is spent to increase the charge of a cell from 95% to 100%.
- Second, a 100% charged cell reliably operates even after several 95% charge restorations, but it needs to be fully charged again after a finite number of 95% charge restorations. This finite number varies from cell to cell.

Based on these observations, the paper defines two types of refresh operations: (1) *full refresh* and (2) *partial refresh*. Full refresh uses the nominal latency and restores the cell charge to 100%, while the latency of partial refresh is only 60% of the nominal value and it restores 95% of the initial charge level. The **key idea** of the paper is to apply a *full refresh* operation **only when necessary** and use *partial refresh* operations at all other times.

(a) Consider a case in which:

- Each row must be refreshed every 64 ms (i.e., the refresh interval is 64 ms).
- Row refresh commands are evenly distributed across the refresh interval. In other words, all rows are refreshed exactly once in *any* given 64 ms time window.
- You are given the following plot, which shows *the distribution of the maximum number of partial refreshes* across all rows of a particular bank.
- We define T as the time that a bank is busy serving the refresh requests in a window of 64ms. if all rows are always fully refreshed (baseline).



How much time does it take (in terms of T) for a bank to refresh all rows within a refresh interval after applying Variable Refresh Latency?

³ Das et al., "VRL-DRAM: Improving DRAM Performance via Variable Refresh Latency." In Proceedings of the 55th Annual Design Automation Conference (DAC), 2018.

(b) You find out that you can relax the refresh interval of the baseline before applying Variable Refresh Latency as follows:

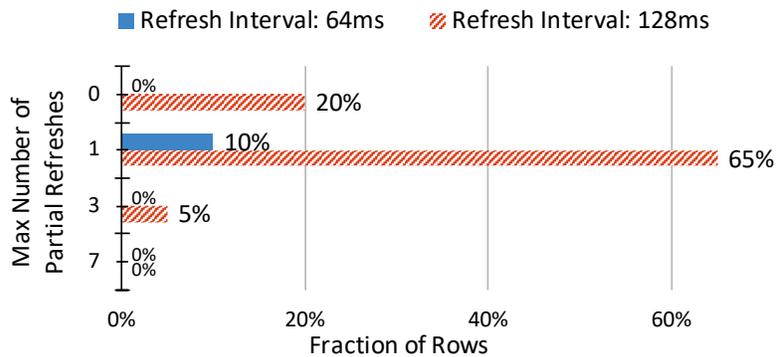
- 90% of the rows are refreshed at every 128ms; 10% of the rows are refreshed at every 64ms.
- Refresh commands are evenly distributed in time.
- All rows are always fully refreshed.
- A row's refresh operation costs $0.2/N$ ms., where N is the number of rows in a bank.

We define *refresh overhead* as the fraction of time that a bank is busy, serving the refresh requests over a very large period of time. Calculate the refresh overhead for the baseline with the relaxed refresh interval.



(c) Consider a case where:

- 90% of the rows are refreshed at every 128ms; 10% of the rows are refreshed at every 64ms.
- Refresh commands are evenly distributed in time.
- You are given the following plot, which shows *the distribution of the maximum number of partial refreshes* across all rows of a particular bank.
- Fully refreshing a row costs $0.2/N$ ms., where N is the number of rows in a bank.
- *Refresh overhead* is defined as the fraction of time that a bank is busy, serving the refresh requests over a very large period of time.



Calculate the refresh overhead and compare it against the baseline configuration (the previous question). How much reduction do you see in the performance overhead of refreshes?