

COMPUTER ARCHITECTURE (263-2210-00L), FALL 2018

HW 4: MEMORY INTERFERENCE AND QoS, EMERGING MEMORY TECHNOLOGIES,
AND PROCESSING-IN-MEMORY

Instructor: Prof. Onur Mutlu

TAs: Mohammed Alser, Can Firtina, Hasan Hassan, Jeremie Kim, Juan Gómez Luna,
Geraldo Francisco de Oliveira, Minesh Patel, Giray Yaglikci

Assigned: Wednesday, Nov 21, 2018

Due: **Wednesday, Dec 5, 2018**

- **Handin - Critical Paper Reviews (1).** You need to submit your reviews to <https://safari.ethz.ch/review/architecture18/>. Please check your inbox. You should have received an email with the password you can use to login to the paper review system. If you have not received any email, please contact comparch@lists.ethz.ch. In the first page after login, you should click in "Architecture - Fall 2018 Home", and then go to "any submitted paper" to see the list of papers.
- **Handin - Questions (2-8).** Please upload your solution to the Moodle (<https://moodle-app2.let.ethz.ch/>) as a single PDF file. **Please use a typesetting software (e.g., LaTeX) or a word processor (e.g., MS Word, LibreOfficeWriter) to generate your PDF file. Feel free to draw your diagrams either using an appropriate software or by hand, and include the diagrams into your solutions PDF.**

1 Critical Paper Reviews [150 points]

Please read the following handout on how to write critical reviews. We will give out extra credit that is worth 0.5% of your total grade for each good review.

- Lecture slides on guidelines for reviewing papers. Please follow this format.
<https://safari.ethz.ch/architecture/fall2018/lib/exe/fetch.php?media=onur-comparch-f18-how-to-do-the-paper-reviews.pdf>
- Some sample reviews can be found here: <https://safari.ethz.ch/architecture/fall2018/doku.php?id=readings>

(a) Write a one-page critical review for the following paper:

- Y. Cai, S. Ghose, E. F. Haratsch, Y. Luo, and O. Mutlu, "Error Characterization, Mitigation, and Recovery in Flash-Memory-Based Solid-State Drives," Proceedings of the IEEE, 2017.
<https://arxiv.org/pdf/1706.08642.pdf>

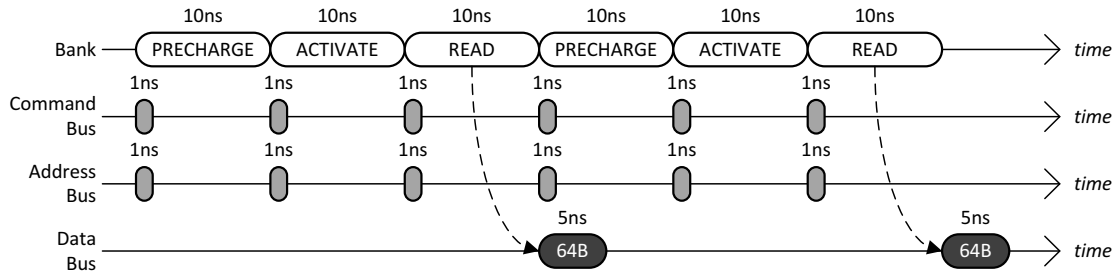
(b) Write a one-page critical review for at least **three** of the following papers:

- O. Mutlu and T. Moscibroda, "Parallelism-Aware Batch Scheduling: Enhancing both Performance and Fairness of Shared DRAM Systems," ISCA 2008. https://people.inf.ethz.ch/omutlu/pub/parbs_isca08.pdf
- E. Ebrahimi, C. J. Lee, O. Mutlu, and Y. N. Patt, "Fairness via Source Throttling: A Configurable and High-Performance Fairness Substrate for Multi-Core Memory Systems," ASPLOS 2010. https://people.inf.ethz.ch/omutlu/pub/fst_asplos10.pdf
- S. P. Muralidhara, L. Subramanian, O. Mutlu, M. Kandemir, and T. Moscibroda, "Reducing Memory Interference in Multicore Systems via Application-Aware Memory Channel Partitioning," MICRO 2011.
<https://people.inf.ethz.ch/omutlu/pub/memory-channel-partitioning-micro11.pdf>

- L. Subramanian, D. Lee, V. Seshadri, H. Rastogi, and O. Mutlu, "BLISS: Balancing Performance, Fairness and Complexity in Memory Access Scheduling," TPDS 2016. https://people.inf.ethz.ch/omutlu/pub/bliss-memory-scheduler_ieee-tpds16.pdf
- A. Boroumand, S. Ghose, Y. Kim, R. Ausavarungnirun, E. Shiu, R. Thakur, D. Kim, A. Kuusela, A. Knies, P. Ranganathan, and O. Mutlu "Google Workloads for Consumer Devices: Mitigating Data Movement Bottlenecks," ASPLOS 2018. https://people.inf.ethz.ch/omutlu/pub/Google-consumer-workloads-data-movement-and-PIM_asplos18.pdf

2 Memory Interference and QoS [100 points]

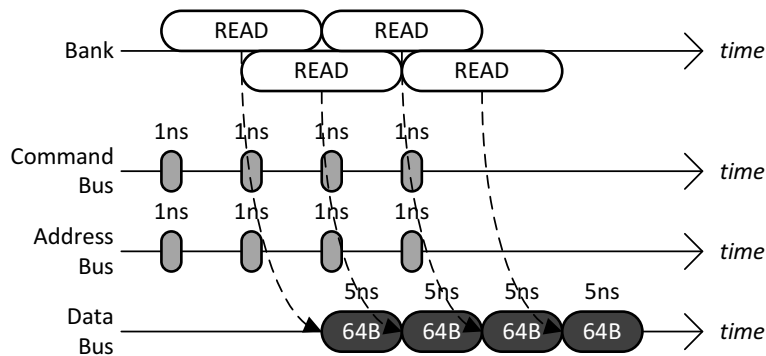
Row-Buffer Conflicts. The following timing diagram shows the operation of a single DRAM channel and a single DRAM bank for two back-to-back reads that conflict in the row-buffer. Immediately after the bank has been busy for 10ns with a READ, data starts to be transferred over the data bus for 5ns.



- (a) Given a long sequence of back-to-back reads that always conflict in the row-buffer, what is the data throughput of the main memory system? Please state your answer in **gigabytes/second**.

- (b) To increase the data throughput, the main memory designer is considering adding more DRAM banks to the single DRAM channel. Given a long sequence of back-to-back reads to all banks that always conflict in the row-buffers, what is the minimum number of banks that is required to achieve the maximum data throughput of the main memory system?

Row-Buffer Hits. The following timing diagram shows the operation of the single DRAM channel and the single DRAM bank for four back-to-back reads that hit in the row-buffer. It is important to note that row-buffer hits to the same DRAM bank are pipelined: while each READ keeps the DRAM bank busy for 10ns, up to at most **half** of this latency (5ns) can be overlapped with another read that hits in the row-buffer. (Note that this is different from Lab 6 where we unrealistically assumed that row-buffer hits are non-pipelined.)



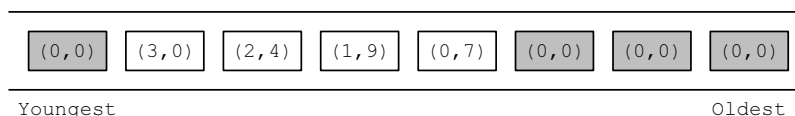
- (c) Given a long sequence of back-to-back reads that always hits in the row-buffer, what is the data throughput of the main memory system? Please state your answer in **gigabytes/second**.

- (d) When the maximum data throughput is achieved for a main memory system that has a single DRAM channel and a single DRAM bank, what is the bottleneck that prevents the data throughput from becoming even larger? **Circle** all that apply.

BANK COMMAND BUS ADDRESS BUS DATA BUS

Memory Scheduling Policies. The diagram below shows the memory controller's *request queue* at time 0. The shaded rectangles are read requests generated by thread T_0 , whereas the unshaded rectangles are read requests generated by thread T_1 . Within each rectangle, there is a pair of numbers that denotes the request's (*BankAddress, RowAddress*). Assume that the memory system has a **single** DRAM channel and four DRAM banks. Further assume the following.

- All the row-buffers are **closed** at time 0.
- Both threads start to stall at time 0 because of memory.
- A thread continues to stall until it receives the data for all of its requests.
- Neither thread generates more requests.



(f) For the *FCFS* scheduling policy, calculate the memory stall time of $T0$ and $T1$.

T0:

T1:

(g) For the *RR – FCFS* scheduling policy, calculate the memory stall time of $T0$ and $T1$.

T0:

T1:

- (h) For the *PAR – BS* scheduling policy, calculate the memory stall time of $T0$ and $T1$. Assume that all eight requests are included in the same batch.

T0:

T1:

3 Memory Scheduling [50 points]

In class, we covered "parallelism-aware batch scheduling," which is a memory scheduling algorithm that aims to reduce interference between threads in a multi-core system.

- (a) What benefit does request batching provide in this algorithm?

- (b) How does the algorithm preserve intra-thread bank parallelism?

- (c) If thread ranking was formed in a "random manner" (i.e., threads were assigned a random rank), would each thread's parallelism be preserved? Why or why not? Explain.

4 Memory Request Scheduling [100 points]

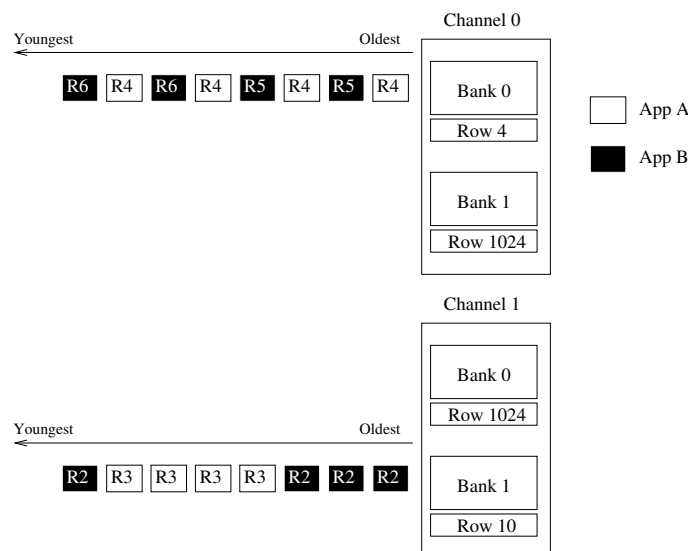
A machine has a DRAM main memory organized as 2 channels, 1 rank and 2 banks/channel. An open row policy is used, i.e., a row is retained in the row-buffer after an access until an access to another row is made. The following commands (defined as we discussed in class) can be issued to DRAM with the given latencies:

- ACTIVATE: 15 ns
- PRECHARGE: 15 ns
- READ/WRITE: 15 ns

Assume the bus latency is 0 cycles.

- (a) Two applications A and B are run on the machine. The following is a snapshot of the request buffers at time t_0 . Requests are tagged with the index of the row they are destined to. Additionally, requests of applications A and B are indicated with different colors. Row 4 is initially open in bank 0 of channel 0 and row 10 is initially open in bank 1 of channel 1.

Each application is stalled until all of its memory requests are serviced and does not generate any more requests.

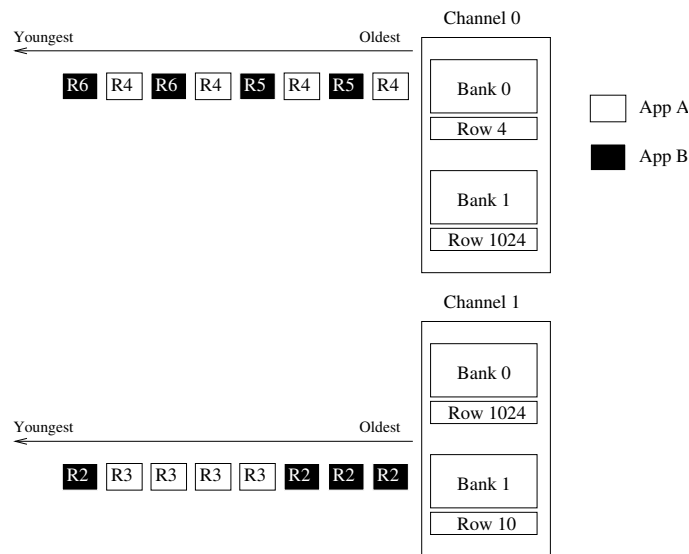


What is the stall time of application A using an FR-FCFS scheduling policy?

What is the stall time of application B using an FR-FCFS scheduling policy?

(b) We studied Parallelism Aware Batch Scheduling (PAR-BS) in class. We will use a PAR-BS-like scheduler that we will call X. This scheduler operates as follows:

- The scheduler forms request batches consisting of the 4 oldest requests from each application at each bank.
- At each bank, the scheduler ranks applications based on the number of requests they have outstanding at that bank. Applications with a smaller number of requests outstanding are assigned a higher rank.
- The scheduler always ranks the application with the oldest request higher in the event of a tie between applications.
- The scheduler prioritizes the requests of applications based on this ranking (higher ranked applications' requests are prioritized over lower ranked applications' requests).
- The scheduler repeats the above steps once all requests in a batch are serviced at all banks.



For the same request buffer state as in Part (a) (replicated above for your benefit):
 What is the stall time of application A using this scheduler?

What is the stall time of application B using this scheduler?

The PAR-BS scheduler we studied in lecture provided better system performance than the FR-FCFS scheduler. Is this true for X also? (i.e., does X provide better system performance than FR-FCFS?)

Circle one:

YES

NO

Explain why or why not. Provide the fundamental reason why X does or does not improve system performance over an FR-FCFS scheduler.

(c) Can you design a better memory scheduler (i.e., one that provides higher system performance) than X?

Circle one:

YES

NO

If yes, answer the questions below. What modifications would you make to scheduler X to design this better scheduler Y? Explain clearly.

What is the stall time of application A using this scheduler Y?

What is the stall time of application B using this scheduler Y?

- (d) Consider a simple channel partitioning scheme where application A's data is mapped to channel 0 and application B's data is mapped to channel 1. When data is mapped to a different channel, only the channel number changes; the bank number does not change. For instance, requests of application A that were mapped to bank 1 of channel 1 would now be mapped to bank 1 of channel 0. What is the stall time of application A using this channel partitioning mechanism and an FR-FCFS memory scheduler?

What is the stall time of application B using this channel partitioning mechanism and an FR-FCFS memory scheduler?

Explain why channel partitioning does better or worse than scheduler Y.

Is channel partitioning and FR-FCFS memory scheduling always strictly better or strictly worse than FR-FCFS memory scheduling alone? Explain.

5 Memory System [80 points]

A machine with a 4 GB DRAM main memory system has 4 channels, 1 rank per channel and 4 banks per rank. The cache block size is 64 bytes.

(a) You are given the following byte addresses and the channel and bank to which they are mapped:

Byte: 0x0000 ⇒ Channel 0, Bank 0
 Byte: 0x0100 ⇒ Channel 0, Bank 0
 Byte: 0x0200 ⇒ Channel 0, Bank 0
 Byte: 0x0400 ⇒ Channel 1, Bank 0
 Byte: 0x0800 ⇒ Channel 2, Bank 0
 Byte: 0x0C00 ⇒ Channel 3, Bank 0
 Byte: 0x1000 ⇒ Channel 0, Bank 1
 Byte: 0x2000 ⇒ Channel 0, Bank 2
 Byte: 0x3000 ⇒ Channel 0, Bank 3

Determine which bits of the address are used for each of the following address components. Assume row bits are higher order than column bits:

- Byte on bus
Addr [2 : 0]
- Channel bits (channel bits are contiguous)
Addr [_____ : _____]
- Bank bits (bank bits are contiguous)
Addr [_____ : _____]
- Column bits (column bits are contiguous)
Addr [_____ : _____]
- Row bits (row bits are contiguous)
Addr [_____ : _____]

(b) Two applications App 1 and App 2 share this memory system (using the address mapping scheme you determined in part (a)). The memory scheduling policy employed is FR-FCFS. The following requests are queued at the memory controller request buffer at time t . Assume the first request (A) is the oldest and the last one ($A + 15$) is the youngest.

A B A + 1 A + 2 A + 3 B + 10 A + 4 B + 12 A + 5 A + 6 A + 7
 A + 8 A + 9 A + 10 + 11 A + 12 A + 13 A + 14 A + 15

These are cache block addresses, not byte addresses. Note that requests to $A + x$ are from App 1, while requests to $B + x$ are from App 2. Addresses A and B are row-aligned (i.e., they are at the start of a row) and are at the same bank but are in different rows.

Assuming row-buffer hits take T time units to service and row-buffer conflicts/misses take $2T$ time units to service, what is the slowdown (compared to when run alone on the same system) of

- App 1?

- App 2?

(c) Which application slows down more?

Why?

(d) In class, we discussed memory channel partitioning and memory request scheduling as two solutions to mitigate interference and application slowdowns in multicore systems. Propose another solution to reduce the slowdown of the more-slowed-down application, without increasing the slowdown of the other application? Be concrete.

6 Emerging Memory Technologies [80 points]

Computer scientists at ETH developed a new memory technology, ETH-RAM, which is non-volatile. The access latency of ETH-RAM is close to that of DRAM while it provides higher density compared to the latest DRAM technologies. ETH-RAM has one shortcoming, however: it has limited endurance, i.e., a memory cell stops functioning after 10^6 writes are performed to the cell (known as cell wear-out).

A bright ETH student has built a computer system using 1 GB of ETH-RAM as main memory. ETH-RAM exploits a perfect wear-leveling mechanism, i.e., a mechanism that equally distributes the writes over all of the cells of the main memory.

(a) This student is worried about the lifetime of the computer system she has built. She executes a test program that runs special instructions to bypass the cache hierarchy and repeatedly writes data into different words until **all** the ETH-RAM cells are worn-out (stop functioning) and the system becomes useless. The student's measurements show that ETH-RAM stops functioning (i.e., all its cells are worn-out) in one year (365 days). Assume the following:

- The processor is in-order and there is no memory-level parallelism.
- It takes 5 ns to send a memory request from the processor to the memory controller and it takes 28 ns to send the request from the memory controller to ETH-RAM.
- ETH-RAM is word-addressable. Thus, each write request writes 4 bytes to memory.

What is the write latency of ETH-RAM? Show your work.

(b) ETH-RAM works in the multi-level cell (MLC) mode in which each memory cell stores 2 bits. The student decides to improve the lifetime of ETH-RAM cells by using the single-level cell (SLC) mode. When ETH-RAM is used in SLC mode, the lifetime of each cell improves by a factor of 10 and the write latency decreases by 70%. What is the lifetime of the system using the SLC mode, if we repeat the experiment in part (a), with everything else remaining the same in the system? Show your work.

7 In-DRAM Bitmap Indices [80 points]

Recall that in class we discussed *Ambit*, which is a DRAM design that can greatly accelerate Bulk Bitwise Operations by providing the ability to perform bitwise AND/OR of two rows in a subarray.

One real-world application that can benefit from *Ambit*'s in-DRAM bulk bitwise operations is the database *bitmap index*, as we also discussed in the lecture. By using bitmap indices, we want to run the following query on a database that keeps track of user actions: "How many unique users were active every week for the past w weeks?" Every week, each user is represented by a single bit. If the user was active a given week, the corresponding bit is set to 1. The total number of users is u .

We assume the bits corresponding to one week are all in the same row. If u is greater than the total number of bits in one row (the row size is 8 kilobytes), more rows in different subarrays are used for the same week. We assume that all weeks corresponding to the users in one subarray fit in that subarray.

We would like to compare two possible implementations of the database query:

- *CPU-based implementation*: This implementation reads the bits of all u users for the w weeks. For each user, it **ands** the bits corresponding to the past w weeks. Then, it performs a bit-count operation to compute the final result.

Since this operation is very memory-bound, we simplify the estimation of the execution time as the time needed to read all bits for the u users in the last w weeks. The memory bandwidth that the CPU can exploit is X bytes/s.

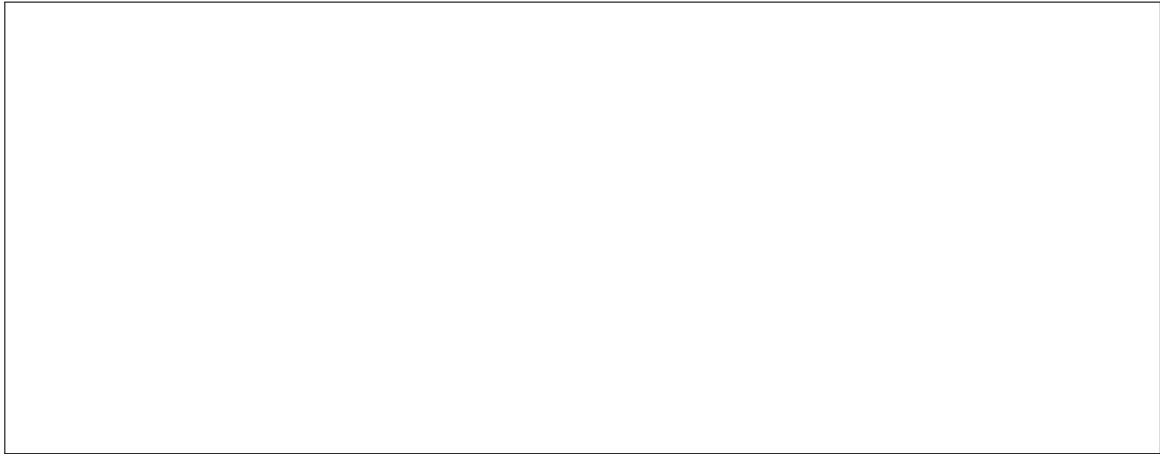
- *Ambit-based implementation*: This implementation takes advantage of bulk **and** operations of *Ambit*. In each subarray, we reserve one *Accumulation* row and one *Operand* row (besides the control rows that are needed for the regular operation of *Ambit*). Initially, all bits in the *Accumulation* row are set to 1. Any row can be moved to the *Operand* row by using *RowClone* (recall that *RowClone* is a mechanism that enables very fast copying of a row to another row in the same subarray). t_{rc} and t_{and} are the latencies (in seconds) of *RowClone*'s copy and *Ambit*'s **and** respectively.

Since *Ambit* does *not* support bit-count operations inside DRAM, the final bit-count is still executed on the CPU. We consider that the execution time of the bit-count operation is negligible compared to the time needed to read all bits from the *Accumulation* rows by the CPU.

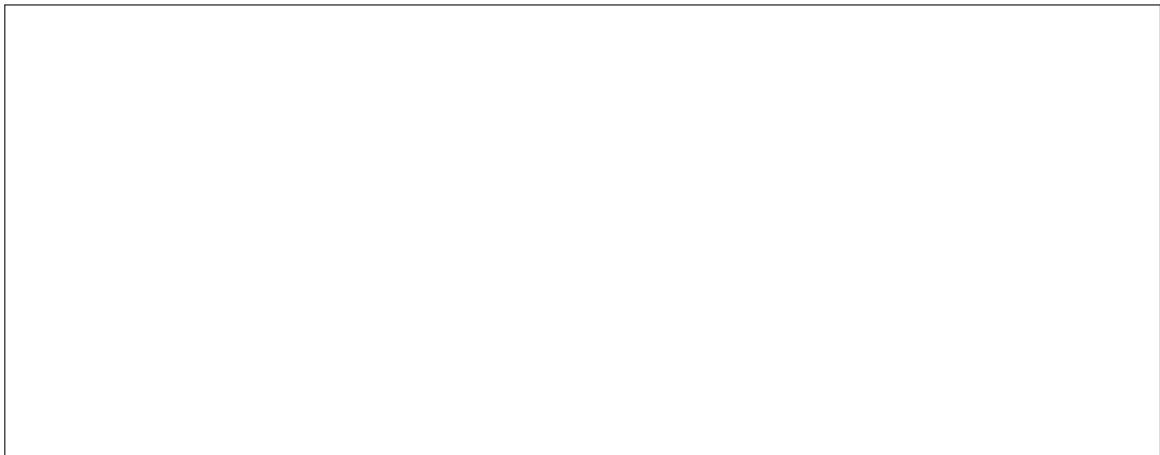
- (a) What is the total number of DRAM rows that are occupied by u users and w weeks?

- (b) What is the throughput in users/second of the *Ambit*-based implementation?

(c) What is the throughput in users/second of the CPU implementation?



(d) What is the maximum w for the CPU implementation to be faster than the Ambit-based implementation? Assume u is a multiple of the row size.



8 Caching vs. Processing-in-Memory [80 points]

We are given the following piece of code that makes accesses to integer arrays A and B. The size of each element in both A and B is 4 bytes. The base address of array A is $0x00001000$, and the base address of B is $0x00008000$.

```
movi R1, #0x1000 // Store the base address of A in R1
movi R2, #0x8000 // Store the base address of B in R2
movi R3, #0

Outer_Loop:
  movi R4, #0
  movi R7, #0
  Inner_Loop:
    add R5, R3, R4 // R5 = R3 + R4
    // load 4 bytes from memory address R1+R5
    ld R5, [R1, R5] // R5 = Memory[R1 + R5],
    ld R6, [R2, R4] // R6 = Memory[R2 + R4]
    mul R5, R5, R6 // R5 = R5 * R6
    add R7, R7, R5 // R7 += R5
    inc R4 // R4++
    bne R4, #2, Inner_Loop // If R4 != 2, jump to Inner_Loop

    //store the data of R7 in memory address R1+R3
    st [R1, R3], R7 // Memory[R1 + R3] = R7,
    inc R3 // R3++
    bne R3, #16, Outer_Loop // If R3 != 16, jump to Outer_Loop
```

You are running the above code on a single-core processor. For now, assume that the processor *does not* have caches. Therefore, all load/store instructions access the main memory, which has a fixed 50-cycle latency, for both read and write operations. Assume that all load/store operations are serialized, i.e., the latency of multiple memory requests *cannot* be overlapped. Also assume that the execution time of a non-memory-access instruction is zero (i.e., we ignore its execution time).

- (a) What is the execution time of the above piece of code in cycles?

- (b) Assume that a 128-byte private cache is added to the processor core in the next-generation processor. The cache block size is 8-byte. The cache is direct-mapped. On a hit, the cache services both read and write requests in 5 cycles. On a miss, the main memory is accessed and the access fills an 8-byte cache line in 50 cycles. Assuming that the cache is initially empty, what is the new execution time on this processor with the described cache? Show your work.



- (c) You are not satisfied with the performance after implementing the described cache. To do better, you consider utilizing a processing unit that is available *close to the main memory*. This processing unit can directly interface to the main memory with a *10-cycle* latency, for both read and write operations. How many cycles does it take to execute the same program using the in-memory processing units? (Assume that the in-memory processing unit does not have a cache, and the memory accesses are serialized like in the processor core. The latency of the non-memory-access operations is ignored.)



- (d) Your friend now suggests that, by changing the cache capacity of the single-core processor (in part (b)), she could provide as good performance as the system that utilizes the memory processing unit (in part (c)).

Is she correct? What is the minimum capacity required for the cache of the single-core processor to match the performance of the program running on the memory processing unit?

- (e) What other changes could be made to the cache design to improve the performance of the single-core processor on this program?