

ETH 263-2210-00L COMPUTER ARCHITECTURE, FALL 2022

HW 1: PROCESSING-IN-MEMORY

Instructor: Prof. Onur Mutlu

TAs: Juan Gómez Luna, Mohammad Sadrosadati, Mohammed Alser, Rahul Bera, Nisa Bostanci, João Dinis Ferreira, Can Firtina, Nika Mansouri Ghiasi, Geraldo Francisco De Oliveira Junior, Konstantinos Kanellopoulos, Joël Lindegger, Rakesh Nadig, Ataberk Olgun, Abdullah Giray Yaglikci, Yahya Can Tugrul, Haocong Luo, Banu Cavlak, Aditya Manglik

Given: Friday, October 7, 2022

Due: **Friday, October 21, 2022**

- **Handin - Critical Paper Reviews (1).** You need to submit your reviews to <https://safari.ethz.ch/review/architecture22/>. Please, check your inbox, you should have received an email with the password you should use to login. If you did not receive any email, contact comparch@lists.inf.ethz.ch. In the first page after login, you should click in "Computer Architecture Home", and then go to "any submitted paper" to see the list of papers.
- **Handin - Questions (2-6).** You should upload your answers to the Moodle Platform (<https://moodle-app2.let.ethz.ch/mod/assign/view.php?id=804436>) as a single PDF file.
- If you have any questions regarding this homework, please ask them the Moodle forum (<https://moodle-app2.let.ethz.ch/mod/moodleoverflow/view.php?id=807326>).
- Please note that the handin questions have a hard deadline. However, you can submit your paper reviews till the end of the semester.

1. Critical Paper Reviews [1,000 points]

You will do at least 5 readings for this homework, out of which 3 are tagged as **REQUIRED** papers. You may access them by *simply clicking on the QR codes below or scanning them*.



Required 1



Required 2



Required 3

Write an approximately one-page critical review for the readings (i.e., papers from #1 to #3 **and at least 2** of the remaining papers, from #4 to #22). If you review a paper other than the 5 mandatory papers, you will receive 200 BONUS points on top of 1,000 points you may get from paper reviews (i.e., each additional submission is worth 200 BONUS points with a possibility to get up to 4400 points). Note that you will get **zero** points from the critical paper reviews if you do not submit the required paper reviews (i.e., papers from #1 to #3).

Please read the guideline slides for reviewing papers and watch Prof. Mutlu's guideline video on how to do a critical paper review. We also provide you with sample reviews which you can access using the QR code. 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. **We will give out extra credit that is worth 0.5% of your total course grade for each good review.**



Guideline Slides



Guideline Video



Sample Reviews

1. **(REQUIRED)** Richard Hamming, "You and Your Research," Bell Communications Research Colloquium Seminar, 1986. <https://https://safari.ethz.ch/architecture/fall2022/lib/exe/fetch.php?media=youandyourresearch.pdf>
2. **(REQUIRED)** Onur Mutlu, "Intelligent Architectures for Intelligent Machines," World Conference on Energy Science and Technology (WCEST), 2021.
 - Talk: https://www.youtube.com/watch?v=20Bdcw-ilQY&list=PL5Q2soXY2Zi8D_5MGV6EnXEJHnV2YFBJl&index=61
 - Slide (PPT): <https://people.inf.ethz.ch/omutlu/pub/onur-TUBA-WCEST-Keynote-IntelligentArchitecturesForIntelligentSystems-August-11-2021.pptx>
 - Slide (PDF): <https://people.inf.ethz.ch/omutlu/pub/onur-TUBA-WCEST-Keynote-IntelligentArchitecturesForIntelligentSystems-August-11-2021.pdf>
3. **(REQUIRED)** Seshadri et al., "Ambit: In-Memory Accelerator for Bulk Bitwise Operations Using Commodity DRAM Technology", MICRO, 2017. https://people.inf.ethz.ch/omutlu/pub/ambit-bulk-bitwise-dram_micro17.pdf
4. Ahn et al., "A Scalable Processing-in-Memory Accelerator for Parallel Graph Processing", ISCA 2015, https://people.inf.ethz.ch/omutlu/pub/tesseract-pim-architecture-for-graph-processing_isca15.pdf
5. Mutlu et al., "A Modern Primer on Processing in Memory", Invited Book Chapter in Emerging Computing: From Devices to Systems - Looking Beyond Moore and Von Neumann, Springer https://people.inf.ethz.ch/omutlu/pub/ModernPrimerOnPIM_springer-emerging-computing-bookchapter21.pdf
6. Ahn et al., "PIM-Enabled Instructions: A Low-Overhead, Locality-Aware Processing-in-Memory Architecture", ISCA 2015, https://people.inf.ethz.ch/omutlu/pub/pim-enabled-instructions-for-low-overhead-pim_isca15.pdf
7. Boroumand et al., "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
8. Singh et al., "FPGA-based Near-Memory Acceleration of Modern Data-Intensive Applications", IEEE Micro 2021, <https://arxiv.org/pdf/2106.06433.pdf>
9. Seshadri et al., "RowClone: Fast and Energy-Efficient In-DRAM Bulk Data Copy and Initialization", MICRO 2013, https://people.inf.ethz.ch/omutlu/pub/rowclone_micro13.pdf
10. Seshadri et al., "Gather-Scatter DRAM: In-DRAM Address Translation to Improve the Spatial Locality of Non-unit Strided Accesses", MICRO 2015, https://people.inf.ethz.ch/omutlu/pub/GSDRAM-gather-scatter-dram_micro15.pdf
11. Wang et al., "FIGARO: Improving System Performance via Fine-Grained In-DRAM Data Relocation and Caching", MICRO 2020, https://people.inf.ethz.ch/omutlu/pub/FIGARO-fine-grained-in-DRAM-data-relocation-and-caching_micro20.pdf
12. Hajinazar and Oliveira et al., "SIMDRAM: An End-to-End Framework for Bit-Serial SIMD Computing in DRAM", ASPLOS 2021, https://people.inf.ethz.ch/omutlu/pub/SIMDRAM_asplos21.pdf
13. Lee et al., "Hardware Architecture and Software Stack for PIM Based on Commercial DRAM Technology", ISCA 2021 <https://ieeexplore.ieee.org/document/9499894>
14. Harold Stone, "A Logic-in-Memory Computer", TC 1970 https://safari.ethz.ch/architecture/fall2020/lib/exe/fetch.php?media=stone_logic_in_memory_1970.pdf
15. Boroumand et al., "Google Neural Network Models for Edge Devices: Analyzing and Mitigating Machine Learning Inference Bottlenecks", PACT 2021 https://people.inf.ethz.ch/omutlu/pub/Google-neural-networks-for-edge-devices-Mensa-Framework_pact21.pdf
16. Giannoula et al., "SynCron: Efficient Synchronization Support for Near-Data-Processing Architectures", HPCA 2021 https://people.inf.ethz.ch/omutlu/pub/SynCron-synchronization-for-near-data-processing-systems_hpca21.pdf
17. Ferreira et al., "pLUTo: Enabling Massively Parallel Computation in DRAM via Lookup Tables", MICRO 2022 <https://arxiv.org/pdf/2104.07699.pdf>
18. Mao et al., "GenPIP: In-Memory Acceleration of Genome Analysis via Tight Integration of Basecalling and Read Mapping", MICRO 2022 <https://arxiv.org/pdf/2209.08600.pdf>
19. Oliveira et al., "DAMOV: A new methodology and benchmark suite for evaluating data movement bottlenecks", IEEE Access 2021 <https://arxiv.org/pdf/2105.03725.pdf>
20. Xie et al., "Processing-in-Memory Enabled Graphics Processors for 3D Rendering", HPCA 2017 <https://ieeexplore.ieee.org/document/7920862>
21. Gokhale et al., "Processing in memory: the Terasys massively parallel PIM array", Computer 1995 <https://ieeexplore.ieee.org/document/375174>
22. Imani et al., "FloatPIM: In-Memory Acceleration of Deep Neural Network Training with High Precision", ISCA 2019 <https://ieeexplore.ieee.org/document/8980299>

2. Caching vs. Processing-in-Memory [200 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?

3. Processing-near-Memory [200 points]

You want to accelerate the following two pieces of code from Application 1 (App1) and Application 2 (App2).

```
// App1. Registers are 4-byte wide
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

Loop:
    ld R4, [R1, R3]        // R4 = MEM[R1 + R3]
    ld R5, [R2, R3]        // R5 = MEM[R2 + R3]
    mult R6, R4, #0xF      // R6 = R4 * 0xF
    add R6, R6, R5         // R6 = R6 + R5
    st [R1, R3], R6        // MEM[R1 + R3] = R6

    inc R3                 // R3++
    bne R3, 1000000000, Loop // If R3 != 1000000000, jump to Loop
```

```
// App2. Registers are 4-byte wide
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
movi R4, #0

Loop:
    ld R5, [R1, R3]        // R5 = MEM[R1 + R3]
    ld R6, [R1, R3, #4]    // R6 = MEM[R1 + R3 + #4]
    sub R5, R5, R6         // R5 = R5 - R6
    mult R5, R5, R5        // R5 = R5 * R5
    add R4, R4, R5         // R4 = R4 + R5
    sqrt R4, R4            // R4 = sqrt(R4)
    st [R2, R3], R4        // MEM[R2 + R3] = R4

    inc R3, #2             // R3=R3+2
    bne R3, 1000000000, Loop // If R3 != 1000000000, jump to Loop
```

We make the following assumptions about the baseline CPU where both applications run:

- The CPU is a single-issue in-order processor and all load/store operations are serialized, i.e., the latency of multiple memory requests cannot be overlapped.
- The clock frequency of the CPU is 1 GHz.
- Each memory operation (i.e., `ld`, `st`) takes 100 ns.
- Each simple arithmetic operation (i.e., `add`, `sub`, `mult`, `inc`) and branch operation (i.e., `bne`) takes 1 clock cycle to execute.
- A complex arithmetic operation (i.e., `sqrt`) takes as long as 50 simple arithmetic operations.
- All memory operations (i.e., `ld`, `st`) and arithmetic operations (i.e., `add`, `sub`, `mult`, `inc`, `sqrt`) operate on 4 bytes of data.

- (a) What is the execution time of App1 and App2 when running on the baseline CPU? Show your work.
Hint: Do *not* account for the execution time of the initial `movi` instructions, since it is negligible in comparison to the loops.

App1:

App2:

You have learned in class that Processing-near-Memory (PnM) architectures can accelerate memory-bound workloads, since they provide higher bandwidth and lower latency than conventional processor-centric architectures. You decide to design PnM accelerators for App1 and App2.

The memory device you use is an early-generation 3D-stacked memory with an internal memory bandwidth of only 40 GB/s (i.e., twice the bandwidth of a single-channel 2D DDR4 memory). The 3D-stacked memory allows you to embed compute resources at the base die of the memory cube, called *logic layer*. However, the logic layer imposes several design limitations:

- The area available to build your accelerator in the logic layer of the 3D-stacked memory is **100 mm²**.
- The maximum clock frequency of the accelerator in the logic layer of the 3D-stacked memory is $\frac{1}{10} \times$ that of baseline CPU.
- The accelerator consists of one or more *processing elements*, each of which contains several *functional units*. Table 1 shows the functional units available to build your accelerators.
- A processing element of the accelerator executes the same computation as an entire iteration of the loop (i.e., *all* the instructions in the loop body). The processing element executes an iteration completely before moving to the next iteration. Therefore, if there are N processing elements, N iterations of the loop are executed in parallel.
- The loop iterations are executed on the processing elements as decided by a compiler, which unrolls the loop and preassigns iterations to the processing elements (i.e., the compiler schedules the iterations statically).
- Each functional unit of a processing element of the accelerator can execute one instruction of the loop body at a time. The same functional unit can execute different instructions (e.g., an arithmetic unit is capable of executing **add**, **sub**, **mult**, **inc**) at different times. Two functional units can execute in parallel, if there is no dependence between the operands.

Table 1. Functional units available to build your accelerators.

| Functional Unit | Description | Latency (cycles) | Area (mm ²) |
|---------------------|---|------------------|-------------------------|
| Arithmetic Unit | Arithmetic unit capable of executing add , sub , mult , inc | 1 | 1 |
| Load and Store Unit | Load and store unit. It can issue 1 4-byte ld/st operation per cycle | 1 | 1 |
| Branch Unit | Executes conditional branches (bne) with 100% accuracy | 1 | 1 |
| Square Root Unit | Executes square root (sqr t) operations | 70 | 30 |

You design your accelerator with the maximum possible number of processing elements, in order to be able to run in parallel as many loop iterations as possible. A processing element can have as many functional units as possible, in order to extract from the code as much Instruction Level Parallelism (ILP) as possible.

- (b) What is the area of the configuration of your PnM accelerator that provides the highest performance for each application while fitting in the PnM area budget? Show your work.

App1:

App2:

- (c) Are the PnM accelerators obtained in (b) capable of fully utilizing the memory bandwidth that the 3D-stacked memory provides for App1 and for App2? Show your work.

App1:

App2:

- (d) What is the speedup of the PnM accelerators compared to the execution of App1 and App2 on the baseline CPU? Show your work.

App1:

App2:

4. Processing in Memory: Ambit [400 points]

4.1. In-DRAM Bitmap Indices I [200 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.

4.2. In-DRAM Bitmap Indices II [200 points]

You have been hired to accelerate ETH's student database. After profiling the system for a while, you found out that one of the most executed queries is to *"select the hometown of the students that are from Switzerland and speak German"*. The attributes *hometown*, *country*, and *language* are encoded using a four-byte binary representation. The database has 32768 (2^{15}) entries, and each attribute is stored contiguously in memory. The database management system executes the following query:

```
1 bool position_hometown[entries];
2 for(int i = 0; i < entries; i++){
3     if(students.country[i] == "Switzerland" && students.language[i] == "German"){
4         position_hometown[i] = true;
5     }
6     else{
7         position_hometown[i] = false;
8     }
9 }
```

- (a) You are running the above code on a single-core processor. Assume that:
- Your processor has an 8 MB direct-mapped cache, with a cache line of 64 bytes.
 - A hit in this cache takes one cycle and a miss takes 100 cycles for both load and store operations.
 - All load/store operations are serialized, i.e., the latency of multiple memory requests cannot be overlapped.
 - The starting addresses of *students.country*, *students.language*, and *position_hometown* are 0x05000000, 0x06000000, 0x07000000 respectively.
 - The execution time of a non-memory instruction is zero (i.e., we ignore its execution time).

How many cycles are required to run the query? Show your work.

- (b) 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/XOR of two rows in a sub-array. AMBIT works by issuing back-to-back ACTIVATE (A) and PRECHARGE (P) operations. For example, to compute AND, OR, and XOR operations, AMBIT issues the sequence of commands described in the table below (e.g., $AAP(X, Y)$ represents double row activation of rows X and Y followed by a precharge operation, $AAAP(X, Y, Z)$ represents triple row activation of rows X, Y, and Z followed by a precharge operation).

In those instructions, AMBIT copies the source rows D_i and D_j to auxiliary rows (B_i). Control rows C_i dictate which operation (AND/OR) AMBIT executes. The DRAM rows with dual-contact cells (i.e., rows DCC_i) are used to perform the bitwise NOT operation on the data stored in the row. Basically, copying a source row to DCC_i flips all bits in the source row and stores the result in both the source row and DCC_i . Assume that:

- The DRAM row size is **8 Kbytes**.
- An ACTIVATE command takes 50 cycles to execute.
- A PRECHARGE command takes 20 cycles to execute.
- DRAM has a single memory bank.
- The syntax of an AMBIT operation is: $bbop_ [and/or/xor] \ destination, source_1, source_2$.
- Addresses 0x08000000 and 0x09000000 are used to store partial results.
- The rows at addresses 0x0A000000 and 0x0B000000 store the codes for "Switzerland" and "German", respectively, in each four bytes throughout the entire row.

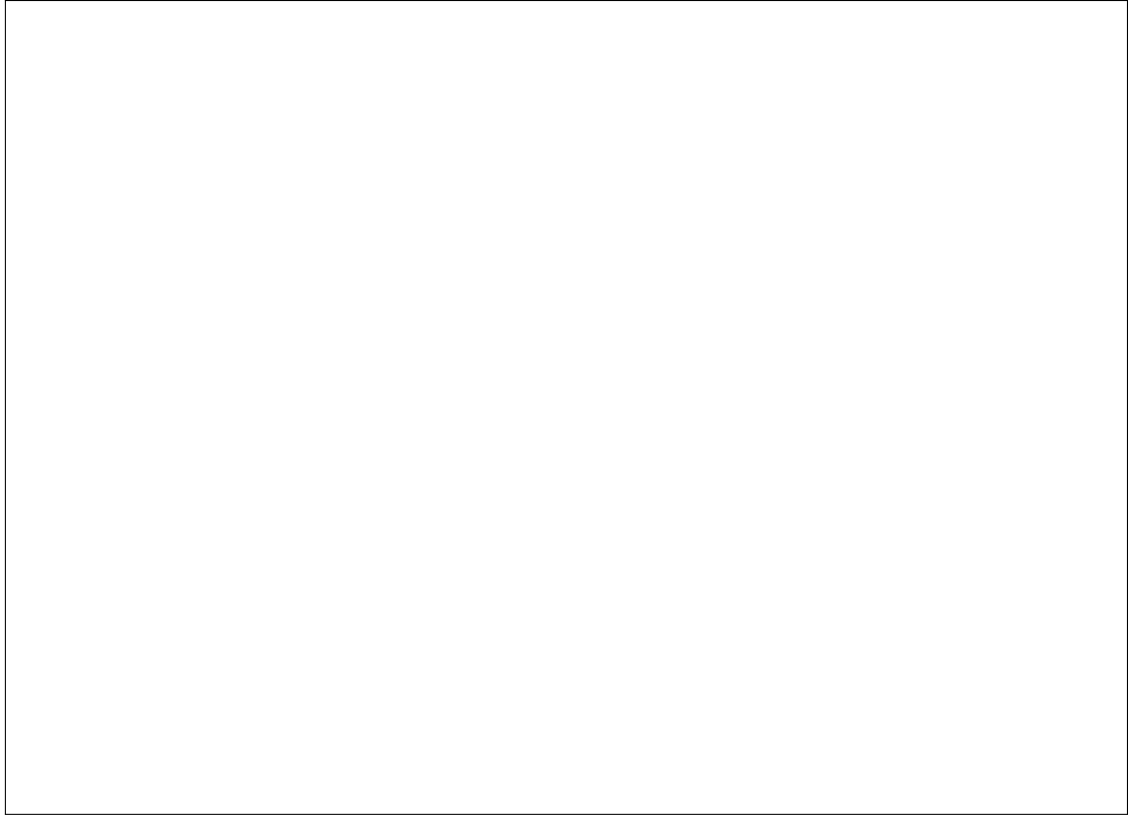
| $D_k = D_i \text{ AND } D_j$ | $D_k = D_i \text{ OR } D_j$ | $D_k = D_i \text{ XOR } D_j$ |
|------------------------------|-----------------------------|------------------------------|
| | | $AAP (D_i, B_0)$ |
| | | $AAP (D_j, B_1)$ |
| | | $AAP (D_i, DCC_0)$ |
| $AAP (D_i, B_0)$ | $AAP (D_i, B_0)$ | $AAP (D_j, DCC_1)$ |
| $AAP (D_j, B_1)$ | $AAP (D_j, B_1)$ | $AAP (C_0, B_2)$ |
| $AAP (C_0, B_2)$ | $AAP (C_1, B_2)$ | $AAAP (B_0, DCC_1, B_2)$ |
| $AAAP (B_0, B_1, B_2)$ | $AAAP (B_0, B_1, B_2)$ | $AAP (C_0, B_2)$ |
| $AAP B_0, D_k$ | $AAP B_0, D_k$ | $AAAP (B_1, DCC_0, B_2)$ |
| | | $AAP (C_1, B_2)$ |
| | | $AAAP (B_0, B_1, B_2)$ |
| | | $AAP (B_0, D_k)$ |

- i) The following code aims to execute the query "select the hometown of the students that are from Switzerland and speak German" in terms of Boolean operations to make use of AMBIT. Fill in the blank boxes such that the algorithm produces the correct result. Show your work.

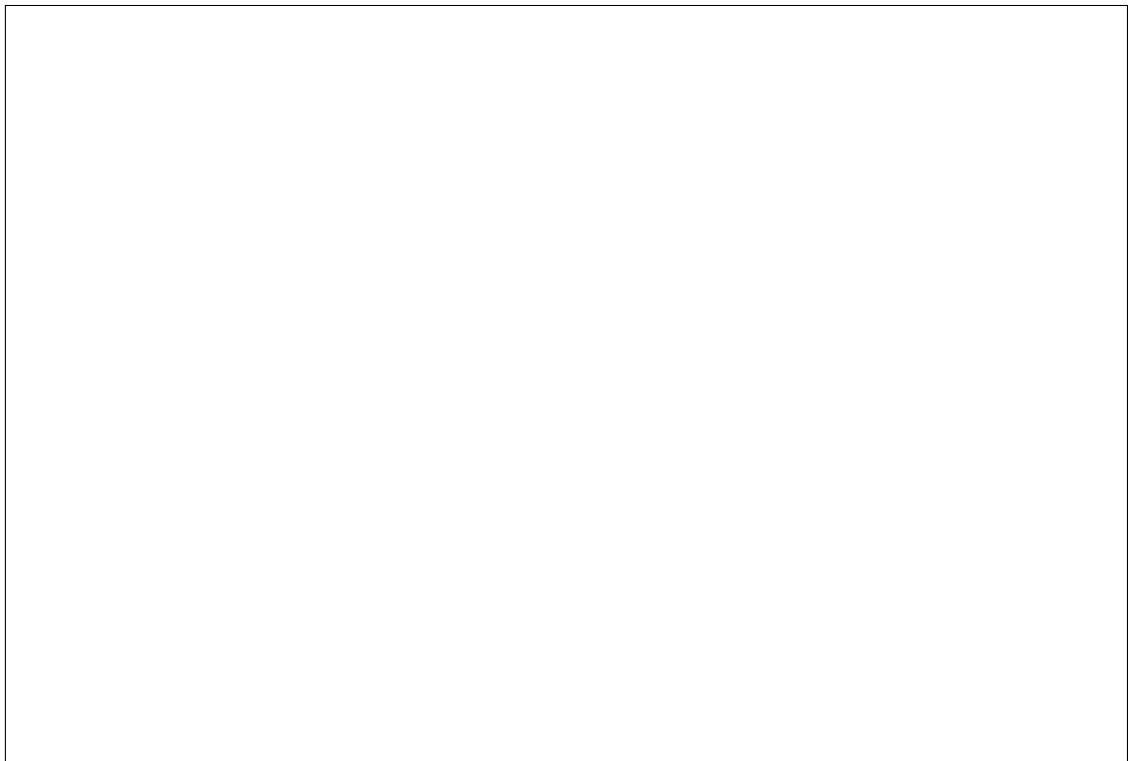
```

1 for(int i = 0; i <  ; i++){
2
3     bbop_ 0x08000000, 0x05000000 + i*8192, 0x0A000000;
4
5     bbop_ 0x09000000, 0x06000000 + i*8192, 0x0B000000;
6
7     bbop_ 0x07000000, 0x08000000, 0x09000000;
8 }

```



- ii) How much speedup does AMBIT provide over the baseline processor when executing the same query? Show your work.



5. In-DRAM Bit Serial Computation [200 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 and NOT of one row. Since Ambit is logically complete, it is possible to implement any other logic gate (e.g., XOR). To be able to implement arithmetic operations, bit shifting is also necessary. There is no way of shifting bits in DRAM with a conventional layout, but it can be done with a bit-serial layout, as Figure 1 shows. With such a layout, it is possible to perform bit-serial arithmetic computations in Ambit.

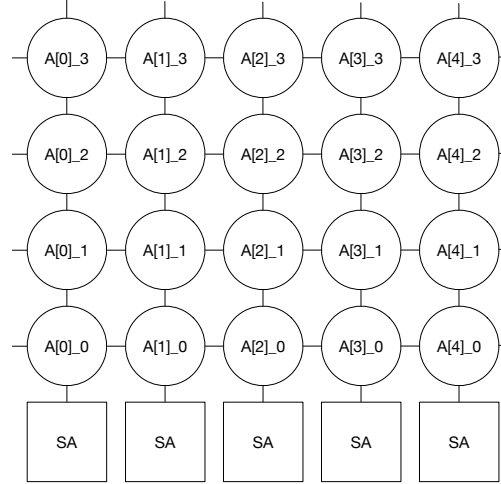


Figure 1. In-DRAM bit-serial layout for array A, which contains five 4-bit elements. DRAM cells in the same bitline contain the bits of an array element: $A[i]_j$ represents bit j of element i .

We want to evaluate the potential performance benefits of using Ambit for arithmetic computations by implementing a simple workload, the element-wise addition of two arrays. Listing 2 shows a sequential code for the addition of two input arrays A and B into output array C.

Listing 1. Sequential CPU implementation of element-wise addition of arrays A and B.

```

1 for(int i = 0; i < num_elements; i++){
2     C[i] = A[i] + B[i];
3 }

```

We compare two possible implementations of the element-wise addition of two arrays: a CPU-based and an Ambit-based implementation. We make two assumptions. First, we use the most favorable layout for each implementation (i.e., conventional layout for CPU, and bit-serial layout for Ambit). Second, both implementations can operate on array elements of any size (i.e., bits/element):

- *CPU-based implementation:* This implementation reads elements of A and B from memory, adds them, and writes the resulting elements of C into memory. Since the computation is simple and regular, we can use a simple analytical performance model for the execution time of the CPU-based implementation: $t_{cpu} = K \times num_operations + \frac{num_bytes}{M}$, where K represents the cost per arithmetic operation and M is the DRAM bandwidth. Note: $num_operations$ should include only the operations for the array addition.
- *Ambit-based implementation:* This implementation assumes a bit serial layout for arrays A, B, and C. It performs additions in a bit serial manner, which only requires XOR, AND, and OR operations, as you can see in the 1-bit full adder in Figure 2.

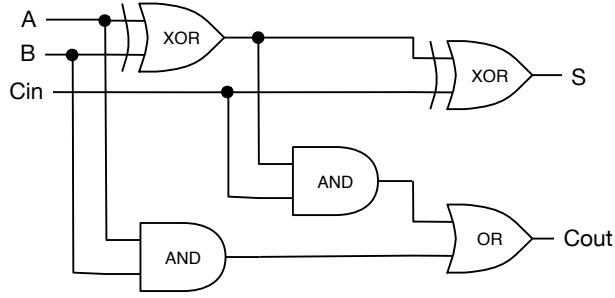


Figure 2. 1-bit full adder.

Ambit implements these operations by issuing back-to-back ACTIVATE (A) and PRECHARGE (P) operations. For example, to compute AND, OR, and XOR operations, Ambit issues the sequence of commands described in Table 2, where $AAAP(X, Y)$ represents double row activation of rows X and Y followed by a precharge operation, and $AAAP(X, Y, Z)$ represents triple row activation of rows X, Y, and Z followed by a precharge operation.

In those instructions, Ambit copies the source rows D_i and D_j to auxiliary rows (B_i). Control rows C_i dictate which operation (AND/OR) Ambit executes. The DRAM rows with dual-contact cells (i.e., rows DCC_i) are used to perform the bitwise NOT operation on the data stored in the row. Basically, the NOT operation copies a source row to DCC_i , flips all bits of the row, and stores the result in both the source row and DCC_i . Assume that:

- The DRAM row size is 8 Kbytes.
- An ACTIVATE command takes 20ns to execute.
- A PRECHARGE command takes 10ns to execute.
- DRAM has a single memory bank.
- The syntax of an Ambit operation is: *bbop_[and/or/xor] destination, source_1, source_2*.
- The rows at addresses 0x00700000, 0x00800000, and 0x00900000 are used to store partial results. Initially, they contain all zeroes.
- The rows at addresses 0x00A00000, 0x00B00000, and 0x00C00000 store arrays A, B, and C, respectively.
- These are all byte addresses. All these rows belong to the same DRAM subarray.

Table 2. Sequences of ACTIVATE and PRECHARGE operations for the execution of Ambit's AND, OR, and XOR.

| $D_k = D_i$ AND D_j | $D_k = D_i$ OR D_j | $D_k = D_i$ XOR D_j |
|------------------------------|-----------------------------|------------------------------|
| | | AAAP (D_i, B_0) |
| | | AAAP (D_j, B_1) |
| | | AAAP (D_i, DCC_0) |
| | | AAAP (D_j, DCC_1) |
| AAAP (D_i, B_0) | AAAP (D_i, B_0) | AAAP (C_0, B_2) |
| AAAP (D_j, B_1) | AAAP (D_j, B_1) | AAAP (C_1, B_2) |
| AAAP (C_0, B_2) | AAAP (C_1, B_2) | AAAP (B_0, DCC_1, B_2) |
| AAAP (B_0, B_1, B_2) | AAAP (B_0, B_1, B_2) | AAAP (C_0, B_2) |
| AAAP B_0, D_k | AAAP B_0, D_k | AAAP (B_1, DCC_0, B_2) |
| | | AAAP (C_1, B_2) |
| | | AAAP (B_0, B_1, B_2) |
| | | AAAP (B_0, D_k) |

- (a) For the CPU-based implementation, you want to obtain K and M . To this end, you run two experiments. In the first experiment, you run your CPU code for the element-wise array addition for 65,536 4-bit elements and measure $t_{cpu} = 100$ us. In the second experiment, you run the STREAM-Copy benchmark for 102,400 4-bit elements and measure $t_{cpu} = 10$ us. The STREAM-Copy benchmark simply copies the contents of one input array **A** to an output array **B**. What are the values of K and M ?

- (b) Write the code for the Ambit-based implementation of the element-wise addition of arrays **A** and **B**. The resulting array is **C**.

- (c) Compute the maximum throughput (in arithmetic operations per second, OPS) of the Ambit-based implementation as a function of the element size (i.e., bits/element).

- (d) Determine the element size (in bits) for which the CPU-based implementation is faster than the Ambit-based implementation (Note: Use the same array size as in the previous part).

6. Processing-using-Memory [200 points]

One promising trend in the Processing-in-Memory paradigm is Processing-using-Memory (PuM), which exploits the analog operation of memory cells to execute bulk bitwise operations. A pioneering proposal in PuM in DRAM technology is Ambit, which we discussed in class. Ambit provides the ability to perform bitwise AND/OR of two rows in a subarray and NOT of one row. Since Ambit is logically complete, it is possible to implement any other logic gate (e.g., XOR). However, to be able to implement arithmetic operations (e.g., addition), bit shifting is also necessary. There is no way of shifting bits in DRAM with a conventional layout, but there are two possible approaches to modifying DRAM to enable bit shifting.

The first approach uses a bit-serial layout (i.e., it changes the horizontal layout to vertical), as Figure 3(a) shows. With such a layout, it is possible to perform **bit-serial** arithmetic computations inside DRAM. For example, performing an addition in a bit-serial manner only requires XOR, AND, and OR operations, as the 1-bit full adder in Figure 3(b) shows.

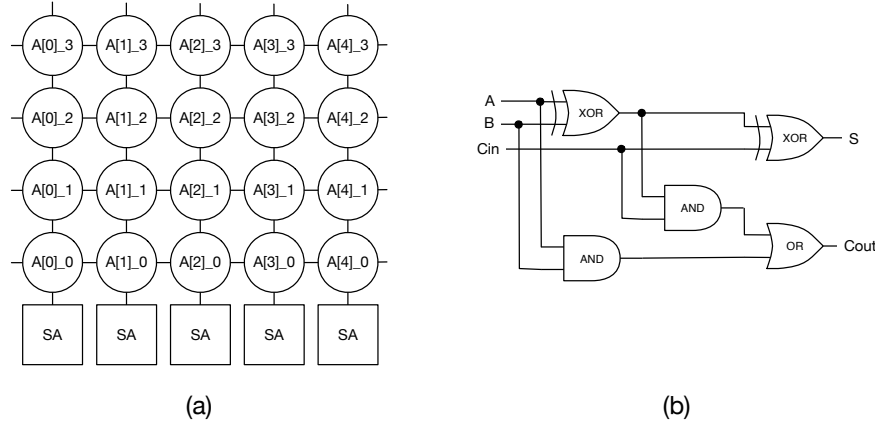


Figure 3. (a) In-DRAM bit-serial layout for array A, which contains five 4-bit elements. DRAM cells in the same bitline contain the bits of an array element: $A[i]_j$ represents bit j of element i . (b) 1-bit full adder.

The second approach uses the conventional horizontal layout, but extends the DRAM subarray with *shifting lines*, which connect each bitline to the previous sense amplifier (SA) to enable left shifting. Figure 4(a) illustrates the second approach, where dashed lines represent the shifting lines. With such shifting lines, it is possible to perform **bit-parallel** arithmetic computations inside DRAM. For example, an addition can be performed in a bit-parallel manner using a parallel adder, such as the Kogge-Stone adder. Figure 4(b) shows an example of an 8-bit Kogge-Stone adder.

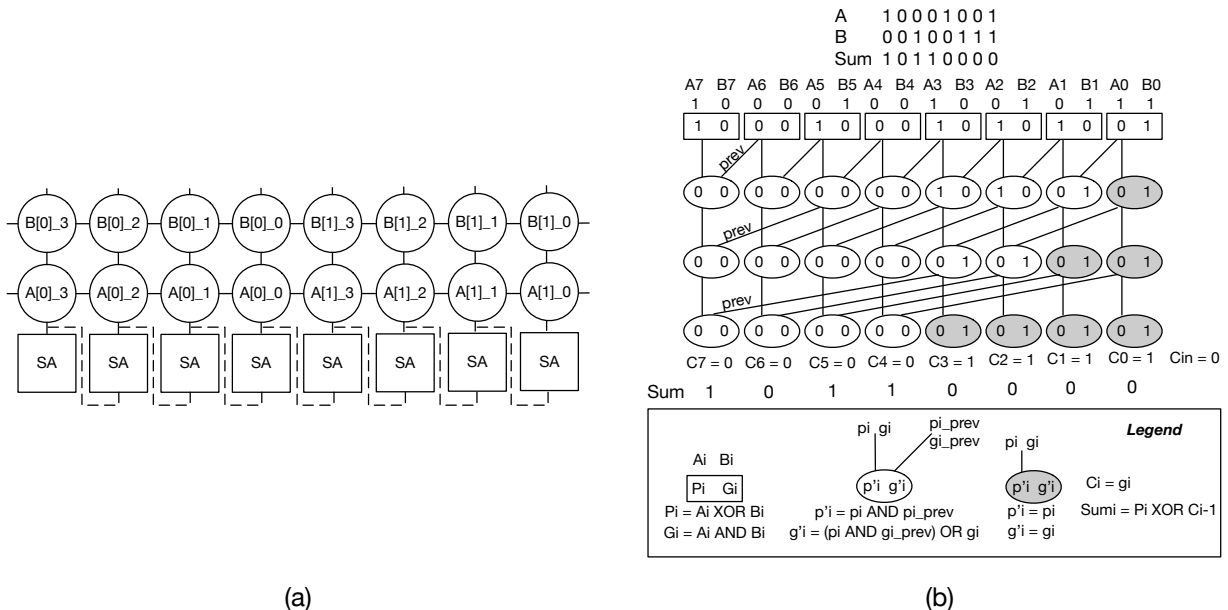


Figure 4. (a) In-DRAM bit-parallel layout for arrays A and B, which contain two 4-bit elements each. DRAM cells in the same bitline contain the same bits of equal-index elements of different arrays. $A[i]_j$ represents bit j of element i . (b) Example of an 8-bit Kogge-Stone adder. A (10001001) and B (00100111) are the two input operands.

We want to compare the potential performance of both approaches to arithmetic computation by implementing a simple workload, the element-wise addition of two arrays. Listing 2 shows a sequential code for the addition of two input arrays A and B into output array C.

Listing 2. Sequential CPU implementation of element-wise addition of arrays A and B.

```
1 for(int i = 0; i < num_elements; i++){
2     C[i] = A[i] + B[i];
3 }
```

As you know from lectures and homeworks, Ambit implements bitwise operations by issuing back-to-back ACTIVATE (A) and PRECHARGE (P) commands. For example, to compute AND, OR, and XOR operations, Ambit issues the sequence of commands described in Figure 5, where AAP(X,Y) represents two consecutive activations of two row addresses X and Y (each of which may correspond to 1, 2, or 3 rows) followed by a precharge operation, and AP(X) represents one activation of row address X followed by a precharge operation.

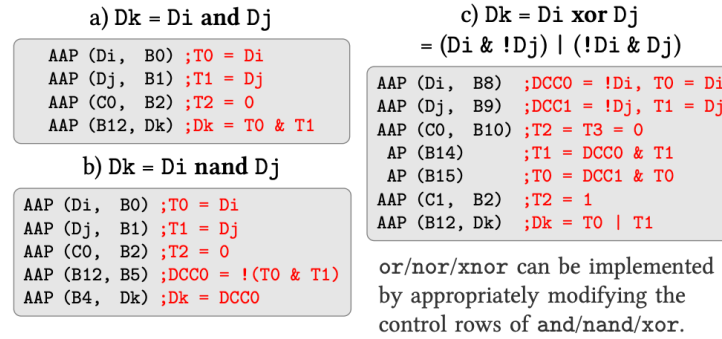


Figure 5. Command sequences for different bitwise operations in Ambit. Notice that AND and OR need the same sequence of commands. Reproduced from Seshadri et al., MICRO 2017.

In those instructions, Ambit copies the source rows Di and Dj to auxiliary row addresses (Bi). Some of the auxiliary row addresses (e.g., B12 in Figure 5) correspond to 3 rows, enabling Triple-Row Activation (TRA), which is the basic operation in Ambit. Control rows Ci dictate which operation (AND/OR) Ambit executes. The DRAM rows with dual-contact cells (i.e., rows DCCi) are used to perform the bitwise NOT operation on the data stored in the row. Basically, the NOT operation copies a source row to DCCi, flips all bits of the row, and stores the result in both the source row and DCCi. Assume that:

- The DRAM row size is 8 Kbytes.
- An ACTIVATE command takes 20ns to execute.
- A PRECHARGE command takes 10ns to execute.
- DRAM has a single memory bank.
- Arrays A, B, and C are properly aligned in both bit-serial and bit-parallel approaches.
- In the bit-parallel approach, a shift operation by one bit requires one AAP.

- (a) Compute the maximum throughput in terms of addition operations per second (OPS) of the bit-serial approach as a function of the element size (i.e., bits/element).

- (b) Compute the maximum throughput in terms of addition operations per second (OPS) of the bit-parallel approach as a function of the element size (i.e., bits/element). Hint: $\sum_{i=0}^n x^i = \frac{1-x^{n+1}}{1-x}$.

- (c) Determine the element size (in bits) for which one approach (i.e., bit-serial or bit-parallel) is preferred over the other one.